
Flow Matching for Scalable Simulation-Based Inference

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Neural posterior estimation methods based on discrete normalizing flows have
2 become established tools for simulation-based inference (SBI), but scaling them
3 to high-dimensional problems can be challenging. Building on recent advances in
4 generative modeling, we here present flow matching posterior estimation (FMPE),
5 a technique for SBI using continuous normalizing flows. Like diffusion models,
6 and in contrast to discrete flows, flow matching allows for unconstrained archite-
7 ctures, providing enhanced flexibility for complex data modalities. Flow matching,
8 therefore, enables exact density evaluation, fast training, and seamless scalability
9 to large architectures—making it ideal for SBI. We show that FMPE achieves com-
10 petitive performance on an established SBI benchmark, and then demonstrate its
11 improved scalability on a challenging scientific problem: for gravitational-wave in-
12 ference, FMPE outperforms methods based on comparable discrete flows, reducing
13 training time by 30% with substantially improved accuracy. Our work underscores
14 the potential of FMPE to enhance performance in challenging inference scenarios,
15 thereby paving the way for more advanced applications to scientific problems.

16 1 Introduction

17 The ability to readily represent Bayesian posteriors of arbitrary complexity using neural networks
18 would herald a revolution in scientific data analysis. Such networks could be trained using simulated
19 data and used for amortized inference across observations—bringing tractable inference and speed to
20 a myriad of scientific models. Thanks to innovative architectures such as normalizing flows [1, 2],
21 approaches to neural simulation-based inference (SBI) [3] have seen remarkable progress in recent
22 years. Here, we show that modern approaches to deep generative modeling (particularly flow
23 matching) deliver substantial improvements in simplicity, flexibility and scaling when adapted to SBI.

24 The Bayesian approach to data analysis is to compare observations to models via the posterior
25 distribution $p(\theta|x)$. This gives our degree of belief that model parameters θ gave rise to an observation
26 x , and is proportional to the model likelihood $p(x|\theta)$ times the prior $p(\theta)$. One is typically interested
27 in representing the posterior in terms of a collection of samples, however obtaining these through
28 standard likelihood-based algorithms can be challenging for intractable or expensive likelihoods. In
29 such cases, SBI offers an alternative based instead on *data simulations* $x \sim p(x|\theta)$. Combined with
30 deep generative modeling, SBI becomes a powerful paradigm for scientific inference [3]. Neural
31 posterior estimation (NPE) [4–6], for instance, trains a conditional density estimator $q(\theta|x)$ to
32 approximate the posterior, allowing for rapid sampling and density estimation for any x consistent
33 with the training distribution.

34 The NPE density estimator $q(\theta|x)$ is commonly taken to be a (discrete) normalizing flow [1, 2], an
35 approach that has brought state-of-the-art performance in challenging problems such as gravitational-
36 wave inference [7]. Naturally, performance hinges on the expressiveness of $q(\theta|x)$. Normalizing

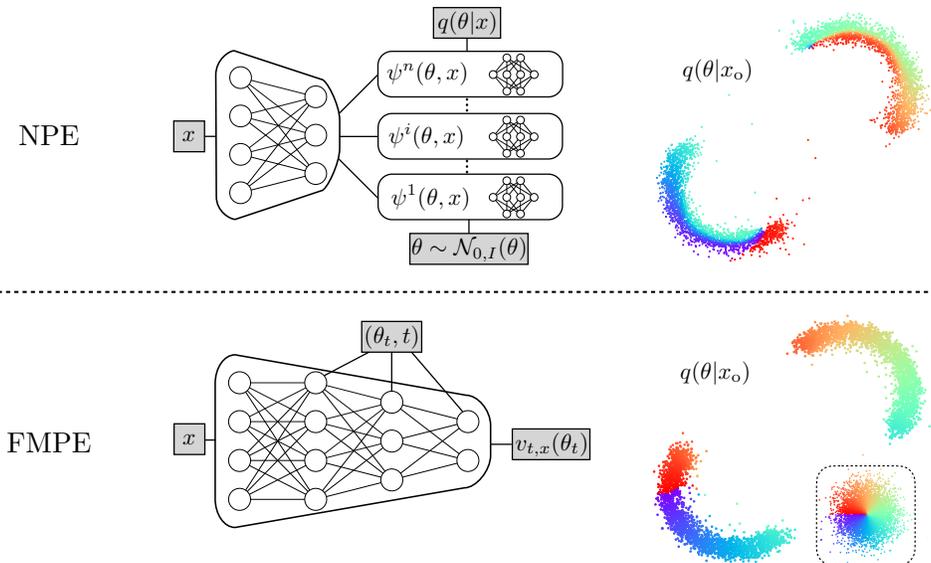


Figure 1: Comparison of network architectures (left) and flow trajectories (right). Discrete flows (NPE, top) require a specialized architecture for the density estimator. Continuous flows (FMPE, bottom) are based on a vector field parametrized with an unconstrained architecture. FMPE uses this additional flexibility to put an enhanced emphasis on the conditioning data x , which in the SBI context is typically high dimensional and in a complex domain. Further, the optimal transport path produces simple flow trajectories from the base distribution (inset) to the target.

37 flows transform noise to samples through a discrete sequence of basic transforms. These have been
 38 carefully engineered to be invertible with simple Jacobian determinant, enabling efficient maximum
 39 likelihood training, while at the same time producing expressive $q(\theta|x)$. Although many such discrete
 40 flows are universal density approximators [2], in practice, they can be challenging to scale to very
 41 large networks, which are needed for big-data experiments.

42 Recent studies [8, 9] propose neural posterior score estimation (NPSE), a rather different approach
 43 that models the posterior distribution with score-matching (or diffusion) networks. These techniques
 44 were originally developed for generative modeling [10–12], achieving state-of-the-art results in many
 45 domains, including image generation [13, 14]. Like discrete normalizing flows, diffusion models
 46 transform noise into samples, but with trajectories parametrized by a *continuous* “time” parameter
 47 t . The trajectories solve a stochastic differential equation [15] (SDE) defined in terms of a vector
 48 field v_t , which is trained to match the score of the intermediate distributions p_t . NPSE has several
 49 advantages compared to NPE, including the ability to combine multiple observations at inference
 50 time [9] and, importantly, the freedom to use unconstrained network architectures.

51 We here propose to use flow matching, another recent technique for generative modeling, for Bayesian
 52 inference, an approach we refer to as flow-matching posterior estimation (FMPE). Flow matching is
 53 also based on a vector field v_t and thereby also admits flexible network architectures (Fig. 1). For flow
 54 matching, however, v_t directly defines the velocity field of sample trajectories, which solve ordinary
 55 differential equations (ODEs) and are deterministic. As a consequence, flow matching allows for
 56 additional freedom in designing non-diffusion paths such as optimal transport, and provides direct
 57 access to the density [16]. These differences are summarized in Tab. 1.

	NPE	NPSE	FMPE (Ours)
Tractable posterior density	Yes	No	Yes
Unconstrained network architecture	No	Yes	Yes
Network passes for sampling	Single	Many	Many

Table 1: Comparison of posterior-estimation methods.

58 Our contributions are as follows:

- 59 • We adapt flow-matching to Bayesian inference, proposing FMPE. In general, the modeling
60 requirements of SBI are different from generative modeling. For the latter, sample quality
61 is critical, i.e., that samples lie in the support of a complex distribution (e.g., images). In
62 contrast, for SBI, $p(\theta|x)$ is typically less complex for fixed x , but x itself can be complex
63 and high-dimensional. We therefore consider pyramid-like architectures from x to v_t , with
64 gated linear units to incorporate (θ, t) dependence, rather than the typical U-Net used for
65 images (Fig. 1). We also propose an alternative t -weighting in the loss, which improves
66 performance in many tasks.
- 67 • Under certain regularity assumptions, we prove an upper bound on the KL divergence
68 between the model and target posterior. This implies that estimated posteriors are mass-
69 covering, i.e., that their support includes all θ consistent with observed x , which is highly
70 desirable for scientific applications [17].
- 71 • We perform a number of experiments to investigate the performance of FMPE. First, on
72 an established suite of SBI benchmarks, we show that FMPE performs comparably—or
73 better—than NPE across most tasks, and in particular exhibits mass-covering posteriors in
74 all cases (Sec. 4). We then turn to the challenging problem of gravitational-wave inference
75 (Sec. 5), where we show that FMPE substantially outperforms an NPE baseline in terms of
76 training time, posterior accuracy, and scaling to larger networks.

77 2 Preliminaries

78 **Normalizing flows.** A normalizing flow [1, 2] defines a probability distribution $q(\theta|x)$ over pa-
79 rameters $\theta \in \mathbb{R}^n$ in terms of an invertible mapping $\psi_x : \mathbb{R}^n \rightarrow \mathbb{R}^n$ from a simple base distribution
80 $q_0(\theta)$,

$$q(\theta|x) = (\psi_x)_* q_0(\theta) = q_0(\psi_x^{-1}(\theta)) \det \left| \frac{\partial \psi_x^{-1}(\theta)}{\partial \theta} \right|, \quad (1)$$

81 where $(\cdot)_*$ denotes the pushforward operator, and for generality we have conditioned on additional
82 context $x \in \mathbb{R}^m$. Unless otherwise specified, a normalizing flow refers to a *discrete* flow, where ψ_x
83 is given by a composition of simpler mappings with triangular Jacobians, interspersed with shuffling
84 of the θ . This construction results in expressive $q(\theta|x)$ and also efficient density evaluation [2].

85 **Continuous normalizing flows.** A continuous flow [18] also maps from base to target distribution,
86 but is parametrized by a continuous “time” $t \in [0, 1]$, where $q_0(\theta|x) = q_0(\theta)$ and $q_1(\theta|x) = q(\theta|x)$.
87 For each t , the flow is defined by a vector field $v_{t,x}$ on the sample space. This corresponds to the
88 velocity of the sample trajectories,

$$\frac{d}{dt} \psi_{t,x}(\theta) = v_{t,x}(\psi_{t,x}(\theta)), \quad \psi_{0,x}(\theta) = \theta. \quad (2)$$

89 We obtain the trajectories $\theta_t \equiv \psi_{t,x}(\theta)$ by integrating this ODE. The final density is given by

$$q(\theta|x) = (\psi_{1,x})_* q_0(\theta) = q_0(\theta) \exp \left(- \int_0^1 \operatorname{div} v_{t,x}(\theta_t) dt \right), \quad (3)$$

90 which is obtained by solving the transport equation $\partial_t q_t + \operatorname{div}(q_t v_{t,x}) = 0$.

91 The advantage of the continuous flow is that $v_{t,x}(\theta)$ can be simply specified by a neural network
92 taking $\mathbb{R}^{n+m+1} \rightarrow \mathbb{R}^n$, in which case (2) is referred to as a *neural ODE* [18]. Since the density
93 is tractable via (3), it is in principle possible to train the flow by maximizing the (log-)likelihood.
94 However, this is often not feasible in practice, since both sampling and density estimation require
95 many network passes to numerically solve the ODE (2).

96 **Flow matching.** An alternative training objective for continuous normalizing flows is provided
97 by flow matching [16]. This directly regresses $v_{t,x}$ on a vector field $u_{t,x}$ that generates a target
98 probability path $p_{t,x}$. It has the advantage that training does not require integration of ODEs, however
99 it is not immediately clear how to choose $(u_{t,x}, p_{t,x})$. The key insight of [16] is that, if the path is

100 chosen on a *sample-conditional* basis,¹ then the training objective becomes extremely simple. Indeed,
 101 given a sample-conditional probability path $p_t(\theta|\theta_1)$ and a corresponding vector field $u_t(\theta|\theta_1)$, we
 102 specify the sample-conditional flow matching loss as

$$\mathcal{L}_{\text{SCFM}} = \mathbb{E}_{t \sim \mathcal{U}[0,1], x \sim p(x), \theta_1 \sim p(\theta|x), \theta_t \sim p_t(\theta|\theta_1)} \|v_{t,x}(\theta_t) - u_t(\theta_t|\theta_1)\|^2. \quad (4)$$

103 Remarkably, minimization of this loss is equivalent to regressing $v_{t,x}(\theta)$ on the *marginal* vector
 104 field $u_{t,x}(\theta)$ that generates $p_t(\theta|x)$ [16]. Note that in this expression, the x -dependence of $v_{t,x}(\theta)$ is
 105 picked up via the expectation value, with the sample-conditional vector field independent of x .

106 There exists considerable freedom in choosing a sample-conditional path. Ref. [16] introduces the
 107 family of Gaussian paths

$$p_t(\theta|\theta_1) = \mathcal{N}(\theta|\mu_t(\theta_1), \sigma_t(\theta_1)^2 I_n), \quad (5)$$

108 where the time-dependent means $\mu_t(\theta_1)$ and standard deviations $\sigma_t(\theta_1)$ can be freely specified
 109 (subject to boundary conditions²). For our experiments, we focus on the optimal transport paths
 110 defined by $\mu_t(\theta_1) = t\theta_1$ and $\sigma_t(\theta_1) = 1 - (1 - \sigma_{\min})t$ (also introduced in [16]). The sample-
 111 conditional vector field then has the simple form

$$u_t(\theta|\theta_1) = \frac{\theta_1 - (1 - \sigma_{\min})\theta}{1 - (1 - \sigma_{\min})t}. \quad (6)$$

112 **Neural posterior estimation (NPE).** NPE is an SBI method that directly fits a density estimator
 113 $q(\theta|x)$ (usually a normalizing flow) to the posterior $p(\theta|x)$ [4–6]. NPE trains with the maximum
 114 likelihood objective $\mathcal{L}_{\text{NPE}} = -\mathbb{E}_{p(\theta)p(x|\theta)} \log q(\theta|x)$, using Bayes’ theorem to simplify the expecta-
 115 tion value with $\mathbb{E}_{p(x)p(\theta|x)} \rightarrow \mathbb{E}_{p(\theta)p(x|\theta)}$. During training, \mathcal{L}_{NPE} is estimated based on an empirical
 116 distribution consisting of samples $(\theta, x) \sim p(\theta)p(x|\theta)$. Once trained, NPE can perform inference
 117 for every new observation using $q(\theta|x)$, thereby *amortizing* the computational cost of simulation
 118 and training across all observations. NPE further provides exact density evaluations of $q(\theta|x)$. Both
 119 of these properties are crucial for the physics application in section 5, so we aim to retain these
 120 properties with FMPE.

121 Related work

122 Flow matching [16] has been developed as a technique for generative modeling, and similar techniques
 123 are discussed in [19–21] and extended in [22, 23]. Flow matching encompasses the deterministic ODE
 124 version of diffusion models [10–12] as a special instance. Although to our knowledge flow matching
 125 has not previously been applied to Bayesian inference, score-matching diffusion models have been
 126 proposed for SBI in [8, 9] with impressive results. These studies, however, use stochastic formulations
 127 via SDEs [15] or Langevin steps and are therefore not directly applicable when evaluations of the
 128 posterior density are desired (see Tab. 1). It should be noted that score modeling can also be used
 129 to parameterize continuous normalizing flows via an ODE. Extension of [8, 9] to the deterministic
 130 formulation could thereby be seen as a special case of flow matching. Many of our analyses and the
 131 practical guidance provided in Section 3 therefore also apply to score matching.

132 We here focus on comparisons of FMPE against NPE [4–6], as it best matches the requirements of
 133 the application in section 5. Other SBI methods include approximate Bayesian computation [24–28],
 134 neural likelihood estimation [29–32] and neural ratio estimation [33–38]. Many of these approaches
 135 have sequential versions, where the estimator networks are specifically tuned to a specific observation
 136 x_0 . FMPE has a tractable density, so it is straightforward to apply the sequential NPE [4–6] approaches
 137 to FMPE. In this case, inference is no longer amortized, so we leave this extension to future work.

138 3 Flow matching posterior estimation

139 To apply flow matching to SBI we use Bayes’ theorem to make the usual replacement $\mathbb{E}_{p(x)p(\theta|x)} \rightarrow$
 140 $\mathbb{E}_{p(\theta)p(x|\theta)}$ in the loss function (4), eliminating the intractable expectation values. This gives the
 141 FMPE loss

$$\mathcal{L}_{\text{FMPE}} = \mathbb{E}_{t \sim p(t), \theta_1 \sim p(\theta), x \sim p(x|\theta_1), \theta_t \sim q_{t,x}(\theta_t|\theta_1)} \|v_{t,x}(\theta_t) - u_{t,x}(\theta_t|\theta_1)\|^2, \quad (7)$$

¹We refer to conditioning on θ_1 as *sample-conditioning* to distinguish from conditioning on x .

²The sample-conditional probability path should be chosen to be concentrated around θ_1 at $t = 1$ (within a small region of size σ_{\min}) and to be the base distribution at $t = 0$.

142 which we minimize using empirical risk minimization over samples $(\theta, x) \sim p(\theta)p(x|\theta)$. In other
 143 words, training data is generated by sampling θ from the prior, and then simulating data x correspond-
 144 ing to θ . This is in close analogy to NPE training, but replaces the log likelihood maximization with
 145 the sample-conditional flow matching objective. Note that in this expression we also sample $t \sim p(t)$,
 146 $t \in [0, 1]$ (see Sec. 3.3), which generalizes the uniform distribution in (4). This provides additional
 147 freedom to improve learning in our experiments.

148 3.1 Probability mass coverage

149 As we show in our examples, trained FMPE models $q(\theta|x)$ can achieve excellent results in approxi-
 150 mating the true posterior $p(\theta|x)$. However, it is not generally possible to achieve *exact* agreement
 151 due to limitations in training budget and network capacity. It is therefore important to understand
 152 how inaccuracies manifest. Whereas sample quality is the main criterion for generative modeling, for
 153 scientific applications one is often interested in the overall shape of the distribution. In particular, an
 154 important question is whether $q(\theta|x)$ is *mass-covering*, i.e., whether it contains the full support of
 155 $p(\theta|x)$. This minimizes the risk to falsely rule out possible explanations of the data. It also allows us
 156 to use importance sampling if the likelihood $p(x|\theta)$ of the forward model can be evaluated, which
 157 can be used for precise estimation of the posterior [39, 40].

158 Consider first the mass-covering property for NPE.
 159 NPE directly minimizes the forward KL divergence
 160 $D_{\text{KL}}(p(\theta|x)||q(\theta|x))$, and thereby provides probability-
 161 mass covering results. Therefore, even if NPE is not accu-
 162 rately trained, the estimate $q(\theta|x)$ should cover the entire
 163 support of the posterior $p(\theta|x)$ and the failure to do so can
 164 be observed in the validation loss. As an illustration in an
 165 unconditional setting, we observe that a unimodal Gaus-
 166 sian q fitted to a bimodal target distribution p captures both
 167 modes when using the forward KL divergence $D_{\text{KL}}(p||q)$,
 168 but only a single mode when using the backwards direction
 169 $D_{\text{KL}}(q||p)$ (Fig. 2).

170 For FMPE, we can fit a Gaussian flow-matching model
 171 $q(\theta) = \mathcal{N}(\hat{\mu}, \hat{\sigma}^2)$ to the same bimodal target, in this case,
 172 parametrizing the vector field as

$$v_t(\theta) = \frac{(\sigma_t^2 + (t\hat{\sigma})^2 - \sigma_t)\theta_t + t\hat{\mu} \cdot \sigma_t}{t \cdot (\sigma_t^2 + (t\hat{\sigma})^2)} \quad (8)$$

173 (see Appendix A), we also obtain a mass-covering distribution when fitting the learnable parameters
 174 $(\hat{\mu}, \hat{\sigma})$ via (4). This provides some indication that the flow matching objective induces mass-covering
 175 behavior, and leads us to investigate the more general question of whether the mean squared error
 176 between vector fields u_t and v_t bounds the forward KL divergence. Indeed, the former agrees up to
 177 constant with the sample-conditional loss (4) (see Sec. 2).

178 We denote the flows of u_t, v_t , by ϕ_t, ψ_t , respectively, and we set $q_t = (\psi_t)_*q_0, p_t = (\phi_t)_*q_0$. The
 179 precise question then is whether we can bound $D_{\text{KL}}(p_1||q_1)$ by $\text{MSE}_p(u, v)^\alpha$ for some positive power
 180 α . It was already observed in [41] that this is not true in general, and we provide a simple example to
 181 that effect in Lemma 1 in Appendix B. Indeed, it was found in [41] that to bound the forward KL
 182 divergence we also need to control the Fisher divergence, $\int p_t(d\theta)(\nabla \ln p_t(\theta) - \nabla q_t(\theta))^2$.

183 Here we show instead that a bound can be obtained under sufficiently strong regularity assumptions
 184 on p_0, u_t , and v_t . The following statement is slightly informal, and we refer to the supplement for the
 185 complete version.

186 **Theorem 1.** *Let $p_0 = q_0$ and assume u_t and v_t are two vector fields whose flows satisfy $p_1 = (\phi_1)_*p_0$
 187 and $q_1 = (\psi_1)_*q_0$. Assume that p_0 is square integrable and satisfies $|\nabla \ln p_0(\theta)| \leq c(1 + |\theta|)$
 188 and u_t and v_t have bounded second derivatives. Then there is a constant $C > 0$ such that (for
 189 $\text{MSE}_p(u, v) < 1$)*

$$D_{\text{KL}}(p_1||q_1) \leq C \text{MSE}_p(u, v)^{\frac{1}{2}}. \quad (9)$$

190 *The proof of this result can be found in appendix B. While the regularity assumptions are not*
 191 *guaranteed to hold in practice when v_t is parametrized by a neural net, the theorem nevertheless*

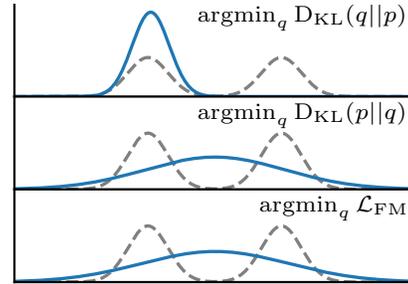


Figure 2: A Gaussian (blue) fitted to a bimodal distribution (gray) with various objectives.

192 gives some indication that the flow-matching objective encourages mass coverage. In Section 4
193 and 5, this is complemented with extensive empirical evidence that flow matching indeed provides
194 mass-covering estimates.

195 We remark that it was shown in [42] that the KL divergence of SDE solutions can be bounded by the
196 MSE of the estimated score function. Thus, the smoothing effect of the noise ensures mass coverage,
197 an aspect that was further studied using the Fokker-Planck equation in [41]. For flow matching,
198 imposing the regularity assumption plays a similar role.

199 3.2 Network architecture

200 Generative diffusion or flow matching models typically operate on complicated and high dimensional
201 data in the θ space (e.g., images with millions of pixels). One typically uses U-Net [43] like
202 architectures, as they provide a natural mapping from θ to a vector field $v(\theta)$ of the same dimension.
203 The dependence on t and an (optional) conditioning vector x is then added on top of this architecture.

204 For SBI, the data x is often associated with a complicated domain, such as image or time series data,
205 whereas parameters θ are typically low dimensional. In this context, it is therefore useful to build the
206 architecture starting as a mapping from x to $v(x)$ and then add conditioning on θ and t . In practice,
207 one can therefore use any established feature extraction architecture for data in the domain of x , and
208 adjust the dimension of the feature vector to $n = \dim(\theta)$. In our experiments, we found that the
209 (t, θ) -conditioning is best achieved using gated linear units [44] to the hidden layers of the network
210 (see also Fig. 1); these are also commonly used for conditioning discrete flows on x .

211 3.3 Re-scaling the time prior

212 The time prior $\mathcal{U}[0, 1]$ in (4) distributes the training capacity uniformly across t . We observed that this
213 is not always optimal in practice, as the complexity of the vector field may depend on t . For FMPE
214 we therefore sample t in (7) from a power-law distribution $p_\alpha(t) \propto t^{1/(1+\alpha)}$, $t \in [0, 1]$, introducing
215 an additional hyperparameter α . This includes the uniform distribution for $\alpha = 0$, but for $\alpha > 0$,
216 assigns greater importance to the vector field for larger values of t . We empirically found this to
217 improve learning for distributions with sharp bounds (e.g., Two Moons in Section 4).

218 4 SBI benchmark

219 We now evaluate FMPE on ten tasks included in the benchmark presented in [45], ranging from
220 simple Gaussian toy models to more challenging SBI problems from epidemiology and ecology, with
221 varying dimensions for parameters ($\dim(\theta) \in [2, 10]$) and observations ($\dim(x) \in [2, 100]$). For
222 each task, we train three separate FMPE models with simulation budgets $N \in \{10^3, 10^4, 10^5\}$. We
223 use a simple network architecture consisting of fully connected residual blocks [46] to parameterize
224 the conditional vector field. For the two tasks with $\dim(x) = 100$ (B-GLM-Raw, SLCP-D), we
225 condition on (t, θ) via gated linear units as described in Section 3.2 (Fig. 8 in Appendix C shows
226 the corresponding performance gain). For the remaining tasks with $\dim(x) \leq 10$ we concatenate
227 (t, θ, x) instead. We reserve 5% of the simulations for validation. See Appendix C for details.

228 For each task and simulation budget, we evaluate the model with the lowest validation loss by
229 comparing $q(\theta|x)$ to the reference posteriors $p(\theta|x)$ provided in [45] for ten different observations x
230 in terms of the C2ST score [47, 48]. This performance metric is computed by training a classifier to
231 discriminate inferred samples $\theta \sim q(\theta|x)$ from reference samples $\theta \sim p(\theta|x)$. The C2ST score is
232 then the test accuracy of this classifier, ranging from 0.5 (best) to 1.0. We observe that FMPE exhibits
233 comparable performance to an NPE baseline model for most tasks and outperforms on several (Fig. 3).
234 In terms of the MMD metric (Fig. 6 in the Appendix), FMPE clearly outperforms NPE (but MMD
235 can be sensitive to its hyperparameters [45]). As NPE is one of the highest ranking methods for many
236 tasks in the benchmark, these results show that FMPE indeed performs competitively with other
237 existing SBI methods. We report an additional baseline for score matching in Fig. 7 in the Appendix.

238 As NPE and FMPE both directly target the posterior with a density estimator (in contrast to most
239 other SBI methods), observed differences can be primarily attributed to their different approaches for
240 density estimation. Interestingly, a great performance improvement of FMPE over NPE is observed
241 for SLCP with a large simulation budget ($N = 10^5$). The SLCP task is specifically designed to have

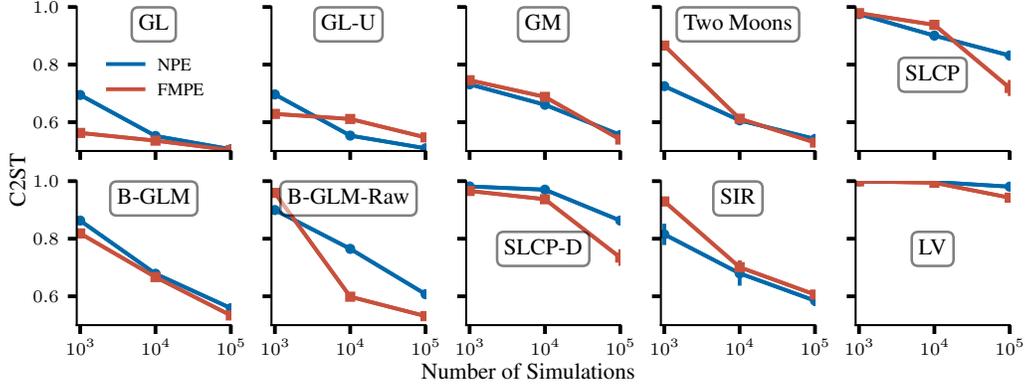


Figure 3: Comparison of FMPE with NPE, a standard SBI method, across 10 benchmark tasks [45].

242 a simple likelihood but a complex posterior, and the FMPE performance underscores the enhanced
 243 flexibility of the FMPE density estimator.

244 Finally, we empirically investigate the mass coverage
 245 suggested by our theoretical analysis in Section 3.1.
 246 We display the density $\log q(\theta|x)$ of the reference
 247 samples $\theta \sim p(\theta|x)$ under our FMPE model q as
 248 a histogram (Fig. 4). All samples $\theta \sim p(\theta|x)$ fall
 249 into the support from $q(\theta|x)$. This becomes appar-
 250 ent when comparing to the density $\log q(\theta|x)$ for
 251 samples $\theta \sim q(\theta|x)$ from q itself. This FMPE re-
 252 sult is therefore mass covering. Note that this does
 253 not necessarily imply conservative posteriors (which
 254 is also not generally true for the forward KL diver-
 255 gence [17, 49, 50]), and some parts of $p(\theta|x)$ may
 256 still be undersampled. Probability mass coverage,
 257 however, implies that no part is entirely missed (com-
 258 pare Fig. 2), even for multimodal distributions such
 259 as Two Moons. Fig. 9 in the Appendix confirms the
 260 mass coverage for the other benchmark tasks.

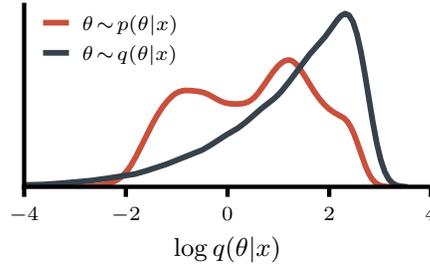


Figure 4: Histogram of FMPE densities $\log q(\theta|x)$ for reference samples $\theta \sim p(\theta|x)$ (Two Moons task, $N = 10^3$). The estimate $q(\theta|x)$ clearly covers $p(\theta|x)$ entirely.

261 5 Gravitational-wave inference

262 5.1 Background

263 Gravitational waves (GWs) are ripples of spacetime predicted by Einstein and produced by cata-
 264 clysmic cosmic events such as the mergers of binary black holes (BBHs). GWs propagate across the
 265 universe to Earth, where the LIGO-Virgo-KAGRA observatories measure faint time-series signals
 266 embedded in noise. To-date, roughly 90 detections of merging black holes and neutron stars have
 267 been made [51], all of which have been characterized using Bayesian inference to compare against
 268 theoretical models.³ These have yielded insights into the origin and evolution of black holes [52],
 269 fundamental properties of matter and gravity [53, 54], and even the expansion rate of the universe [55].

270 Under reasonable assumptions on detector noise, the GW likelihood is tractable,⁴ and inference is
 271 typically performed using tools [56–59] based on Markov chain Monte Carlo [60, 61] or nested

³BBH parameters $\theta \in \mathbb{R}^{15}$ include black-hole masses, spins, and the spacetime location and orientation of the system (see Tab. 4 in the Appendix). We represent x in frequency domain; for two LIGO detectors and complex $f \in [20, 512]$ Hz, $\Delta f = 0.125$ Hz, we have $x \in \mathbb{R}^{15744}$.

⁴Noise is assumed to be stationary and Gaussian, so for frequency-domain data, the GW likelihood $p(x|\theta) = \mathcal{N}(h(\theta)|S_n)(x)$. Here $h(\theta)$ is a theoretical signal model based on Einstein’s theory of general relativity, and S_n is the power spectral density of the detector noise.

272 sampling [62] algorithms. This can take from hours to months, depending on the nature of the event
 273 and the complexity of the signal model, with a typical analysis requiring up to $\sim 10^8$ likelihood
 274 evaluations. The ever-increasing rate of detections means that these analysis times risk becoming a
 275 bottleneck. SBI offers a promising solution for this challenge that has thus been actively studied in
 276 the literature [63–67, 7, 68, 69, 40]. A fully amortized NPE-based method called DINGO recently
 277 achieved accuracies comparable to stochastic samplers with inference times of less than a minute
 278 per event [7]. To achieve accurate results, however, DINGO uses group-equivariant NPE [7, 68]
 279 (GNPE), an NPE extension that integrates known conditional symmetries. GNPE, therefore, does
 280 not provide a tractable density, which is problematic when verifying and correcting inference results
 281 using importance sampling [40].

282 5.2 Experiments

283 We here apply FMPE to GW inference. As a baseline, we train an NPE network with the settings
 284 described in [7] with a few minor changes (see Appendix D).⁵ This uses an embedding network [70]
 285 to compress x to a 128-dimensional feature vector, which is then used to condition a neural spline
 286 flow [71]. The embedding network consists of a learnable linear layer initialized with principal
 287 components of GW simulations followed by a series of dense residual blocks [46]. This architecture
 288 is a powerful feature extractor for GW measurements [7]. As pointed out in Section 3.2, it is
 289 straightforward to reuse such architectures for FMPE, with the following three modifications: (1)
 290 we provide the conditioning on (t, θ) to the network via gated linear units in each hidden layer;
 291 (2) we change the dimension of the final feature vector to the dimension of θ so that the network
 292 parameterizes the conditional vector field $(t, x, \theta) \rightarrow v_{t,x}(\theta)$; (3) we increase the number and width
 293 of the hidden layers to use the capacity freed up by removing the discrete normalizing flow.

294 We train the NPE and FMPE networks with $5 \cdot 10^6$ simulations for 400 epochs using a batch size of
 295 4096 on an A100 GPU. The FMPE network ($1.9 \cdot 10^8$ learnable parameters, training takes ≈ 2 days)
 296 is larger than the NPE network ($1.3 \cdot 10^8$ learnable parameters, training takes ≈ 3 days), but trains
 297 substantially faster. We evaluate both networks on GW150914 [72], the first detected GW. We
 298 generate a reference posterior using the method described in [40]. Fig. 5 compares the inferred
 299 posterior distributions qualitatively and quantitatively in terms of the Jensen-Shannon divergence
 300 (JSD) to the reference.⁶

301 FMPE substantially outperforms NPE in terms of accuracy, with a mean JSD of 0.5 mnat (NPE:
 302 3.6 mnat), and $\max \text{JSD} < 2.0$ mnat, an indistinguishability criterion for GW posteriors [58].
 303 Remarkably, FMPE accuracy is even comparable to GNPE, which leverages physical symmetries
 304 to simplify data. Finally, we find that the Bayesian evidences inferred with NPE ($\log p(x) =$
 305 -7667.958 ± 0.006) and FMPE ($\log p(x) = -7667.969 \pm 0.005$) are consistent within their statistical
 306 uncertainties. A correct evidence is only obtained in importance sampling when the inferred posterior
 307 $q(\theta|x)$ covers the entire posterior $p(\theta|x)$ [40], so this is another indication that FMPE indeed induces
 308 mass-covering posteriors.

309 5.3 Discussion

310 Our results for GW150914 show that FMPE substantially outperforms NPE on this challenging
 311 problem. We believe that this is related to the network structure as follows. The NPE network
 312 allocates roughly two thirds of its parameters to the discrete normalizing flow and only one third
 313 to the embedding network (i.e., the feature extractor for x). Since FMPE parameterizes a much
 314 simpler vector field, it can devote its network capacity to the interpretation of the high-dimensional
 315 $x \in \mathbb{R}^{15744}$, and thereby scales much better to larger networks and achieve much higher accuracy.
 316 Remarkably, the performance is even comparable to GNPE, which involves a much simpler learning
 317 task with likelihood symmetries integrated by construction. See Appendix D for further details.

⁵Our implementation builds on the public DINGO code from <https://github.com/dingo-gw/dingo>.

⁶We omit the three parameters $\phi_c, \phi_{JL}, \theta_{JN}$ in the evaluation as we use phase marginalization in importance sampling and the reference therefore uses a different basis for these parameters [40]. For GNPE we report the results from [7], which are generated with slightly different data conditioning. Therefore, we do not display the GNPE results in the corner plot, and the JSDs serve only as a rough comparison. The JSD for the t_c parameter is not reported in [7] due to a t_c marginalized reference.

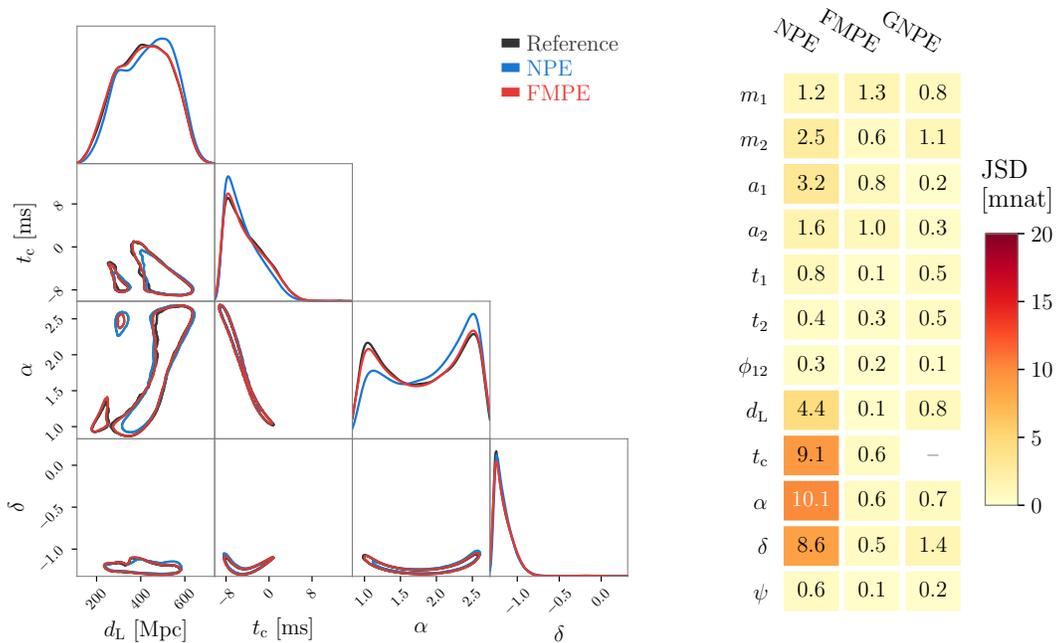


Figure 5: Results for GW150914 [72]. Left: Corner plot showing 1D marginals on the diagonal and 2D 50% credible regions. We display four GW parameters (distance d_L , time of arrival t_c , and sky coordinates α , δ); these represent the least accurate NPE parameters. Right: Deviation between inferred posteriors and the reference, quantified by the Jensen-Shannon divergence (JSD). The FMPE posterior matches the reference more accurately than NPE, and performs similarly to symmetry-enhanced GNPE. (We do not display GNPE results on the left due to different data conditioning settings in available networks.)

318 In future work we plan to carry out a more complete analysis of GW inference using FMPE. Indeed,
 319 GW150914 is a loud event with good data quality, where NPE already performs quite well. DINGO
 320 with GNPE has been validated in a variety of settings [7, 68, 40, 73] including events with a larger
 321 performance gap between NPE and GNPE [68]. Since FMPE (like NPE) does not integrate physical
 322 symmetries, it would likely need further enhancements to fully compete with GNPE. This may require
 323 a symmetry-aware architecture [74], or simply further scaling to larger networks. Nevertheless, our
 324 results demonstrate that FMPE is a promising direction for future research in this field.

325 6 Conclusions

326 We introduced flow matching posterior estimation, a new simulation-based inference technique based
 327 on continuous normalizing flows. In contrast to existing neural posterior estimation methods, it does
 328 not rely on restricted density estimation architectures such as discrete normalizing flows, and instead
 329 parametrizes a distribution in terms of a conditional vector field. This enables more flexible network
 330 architectures and seamless scaling (like score matching), while enabling flexible path specification
 331 and direct access to the posterior density.

332 We evaluated FMPE on a set of 10 benchmark tasks and found competitive or better performance
 333 compared to other simulation-based inference methods. On the challenging task of gravitational-wave
 334 inference, FMPE substantially outperformed comparable discrete flows, producing samples on par
 335 with a method that explicitly leverages symmetries to simplify training. Additionally, flow matching
 336 latent spaces are more naturally structured than those of discrete flows, particularly when using
 337 paths such as optimal transport. Looking forward, it would be interesting to exploit such structure
 338 in designing learning algorithms. This performance and flexibility underscores the capability of
 339 continuous normalizing flows to efficiently solve inverse problems.

References

- 340
- 341 [1] Danilo Rezende and Shakir Mohamed. Variational inference with normalizing flows. In
342 *International Conference on Machine Learning*, pages 1530–1538, 2015. arXiv:1505.05770.
- 343 [2] George Papamakarios, Eric Nalisnick, Danilo Jimenez Rezende, Shakir Mohamed, and Balaji
344 Lakshminarayanan. Normalizing flows for probabilistic modeling and inference. *Journal of*
345 *Machine Learning Research*, 22(57):1–64, 2021. URL: [http://jmlr.org/papers/v22/](http://jmlr.org/papers/v22/19-1028.html)
346 [19-1028.html](http://jmlr.org/papers/v22/19-1028.html).
- 347 [3] Kyle Cranmer, Johann Brehmer, and Gilles Louppe. The frontier of simulation-based inference.
348 *Proc. Nat. Acad. Sci.*, 117(48):30055–30062, 2020. arXiv:1911.01429, doi:10.1073/
349 [pnas.1912789117](https://doi.org/10.1073/pnas.1912789117).
- 350 [4] George Papamakarios and Iain Murray. Fast ε -free inference of simulation models with Bayesian
351 conditional density estimation. In *Advances in neural information processing systems*, 2016.
352 arXiv:1605.06376.
- 353 [5] Jan-Matthis Lueckmann, Pedro J Gonçalves, Giacomo Bassetto, Kaan Öcal, Marcel Non-
354 nenmacher, and Jakob H Macke. Flexible statistical inference for mechanistic models of
355 neural dynamics. In *Proceedings of the 31st International Conference on Neural Information*
356 *Processing Systems*, pages 1289–1299, 2017.
- 357 [6] David Greenberg, Marcel Nonnenmacher, and Jakob Macke. Automatic posterior transformation
358 for likelihood-free inference. In *International Conference on Machine Learning*, pages 2404–
359 2414. PMLR, 2019.
- 360 [7] Maximilian Dax, Stephen R. Green, Jonathan Gair, Jakob H. Macke, Alessandra Buonanno, and
361 Bernhard Schölkopf. Real-Time Gravitational Wave Science with Neural Posterior Estimation.
362 *Phys. Rev. Lett.*, 127(24):241103, 2021. arXiv:2106.12594, doi:10.1103/PhysRevLett.
363 [127.241103](https://doi.org/10.1103/PhysRevLett.127.241103).
- 364 [8] Louis Sharrock, Jack Simons, Song Liu, and Mark Beaumont. Sequential neural score estima-
365 tion: Likelihood-free inference with conditional score based diffusion models. *arXiv preprint*
366 *arXiv:2210.04872*, 2022.
- 367 [9] Tomas Geffner, George Papamakarios, and Andriy Mnih. Score modeling for simulation-based
368 inference. *arXiv preprint arXiv:2209.14249*, 2022.
- 369 [10] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsuper-
370 vised learning using nonequilibrium thermodynamics. In *International Conference on Machine*
371 *Learning*, pages 2256–2265. PMLR, 2015.
- 372 [11] Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data
373 distribution. *Advances in neural information processing systems*, 32, 2019.
- 374 [12] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances*
375 *in Neural Information Processing Systems*, 33:6840–6851, 2020.
- 376 [13] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis.
377 In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors,
378 *Advances in Neural Information Processing Systems*, volume 34, pages 8780–8794. Curran
379 Associates, Inc., 2021. URL: [https://proceedings.neurips.cc/paper_files/paper/](https://proceedings.neurips.cc/paper_files/paper/2021/file/49ad23d1ec9fa4bd8d77d02681df5cfa-Paper.pdf)
380 [2021/file/49ad23d1ec9fa4bd8d77d02681df5cfa-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2021/file/49ad23d1ec9fa4bd8d77d02681df5cfa-Paper.pdf).
- 381 [14] Jonathan Ho, Chitwan Saharia, William Chan, David J. Fleet, Mohammad Norouzi, and Tim
382 Salimans. Cascaded diffusion models for high fidelity image generation. *J. Mach. Learn. Res.*,
383 [23:47:1–47:33](https://doi.org/10.48550/jmlr.2022.23.47.1-47.33), 2022. URL: <http://jmlr.org/papers/v23/21-0635.html>.
- 384 [15] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and
385 Ben Poole. Score-based generative modeling through stochastic differential equations. *arXiv*
386 *preprint arXiv:2011.13456*, 2020.

- 387 [16] Yaron Lipman, Ricky T. Q. Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow
388 matching for generative modeling. *CoRR*, abs/2210.02747, 2022. arXiv:2210.02747, doi:
389 10.48550/arXiv.2210.02747.
- 390 [17] Joeri Hermans, Arnaud Delaunoy, François Rozet, Antoine Wehenkel, and Gilles Louppe.
391 Averting a crisis in simulation-based inference. *arXiv preprint arXiv:2110.06581*, 2021.
- 392 [18] Tian Qi Chen, Yulia Rubanova, Jesse Bettencourt, and David Duvenaud. Neural ordi-
393 nary differential equations. In Samy Bengio, Hanna M. Wallach, Hugo Larochelle,
394 Kristen Grauman, Nicolò Cesa-Bianchi, and Roman Garnett, editors, *Advances in Neu-
395 ral Information Processing Systems 31: Annual Conference on Neural Information
396 Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*,
397 pages 6572–6583, 2018. URL: [https://proceedings.neurips.cc/paper/2018/hash/
398 69386f6bb1dfed68692a24c8686939b9-Abstract.html](https://proceedings.neurips.cc/paper/2018/hash/69386f6bb1dfed68692a24c8686939b9-Abstract.html).
- 399 [19] Michael S Albergo and Eric Vanden-Eijnden. Building normalizing flows with stochastic
400 interpolants. *arXiv preprint arXiv:2209.15571*, 2022.
- 401 [20] Xingchao Liu, Chengyue Gong, and Qiang Liu. Flow straight and fast: Learning to generate
402 and transfer data with rectified flow. *arXiv preprint arXiv:2209.03003*, 2022.
- 403 [21] Kirill Neklyudov, Daniel Severo, and Alireza Makhzani. Action matching: A variational method
404 for learning stochastic dynamics from samples. *arXiv preprint arXiv:2210.06662*, 2022.
- 405 [22] Aram Davtyan, Sepehr Sameni, and Paolo Favaro. Randomized conditional flow matching for
406 video prediction. *arXiv preprint arXiv:2211.14575*, 2022.
- 407 [23] Alexander Tong, Nikolay Malkin, Guillaume Hugué, Yanlei Zhang, Jarrid Rector-Brooks,
408 Kilian Fatras, Guy Wolf, and Yoshua Bengio. Conditional flow matching: Simulation-free
409 dynamic optimal transport. *arXiv preprint arXiv:2302.00482*, 2023.
- 410 [24] Scott A Sisson, Yanan Fan, and Mark A Beaumont. Overview of abc. In *Handbook of*
411 *approximate Bayesian computation*, pages 3–54. Chapman and Hall/CRC, 2018.
- 412 [25] Mark A Beaumont, Wenyang Zhang, and David J Balding. Approximate bayesian computation
413 in population genetics. *Genetics*, 162(4):2025–2035, 2002.
- 414 [26] Mark A Beaumont, Jean-Marie Cornuet, Jean-Michel Marin, and Christian P Robert. Adaptive
415 approximate bayesian computation. *Biometrika*, 96(4):983–990, 2009.
- 416 [27] Michael G. B. Blum and Olivier François. Non-linear regression models for ap-
417 proximate bayesian computation. *Stat. Comput.*, 20(1):63–73, 2010. doi:10.1007/
418 s11222-009-9116-0.
- 419 [28] Dennis Prangle, Paul Fearnhead, Murray P. Cox, Patrick J. Biggs, and Nigel P. French. Semi-
420 automatic selection of summary statistics for abc model choice, 2013. arXiv:1302.5624.
- 421 [29] Simon Wood. Statistical inference for noisy nonlinear ecological dynamic systems. *Nature*,
422 466:1102–4, 08 2010. doi:10.1038/nature09319.
- 423 [30] Christopher C Drovandi, Clara Grazian, Kerrie Mengersen, and Christian Robert. Approximat-
424 ing the likelihood in approximate bayesian computation, 2018. arXiv:1803.06645.
- 425 [31] George Papamakarios, David Sterratt, and Iain Murray. Sequential neural likelihood: Fast
426 likelihood-free inference with autoregressive flows. In *The 22nd International Conference on*
427 *Artificial Intelligence and Statistics*, pages 837–848. PMLR, 2019.
- 428 [32] Jan-Matthis Lueckmann, Giacomo Bassetto, Theofanis Karaletsos, and Jakob H Macke.
429 Likelihood-free inference with emulator networks. In *Symposium on Advances in Approx-
430 imate Bayesian Inference*, pages 32–53. PMLR, 2019.
- 431 [33] Rafael Izbicki, Ann Lee, and Chad Schafer. High-dimensional density ratio estimation with
432 extensions to approximate likelihood computation. In *Artificial intelligence and statistics*, pages
433 420–429. PMLR, 2014.

- 434 [34] Kim Pham, David Nott, and Sanjay Chaudhuri. A note on approximating abc-mcmc using
435 flexible classifiers. *Stat*, 3, 03 2014. doi:10.1002/sta4.56.
- 436 [35] Kyle Cranmer, Juan Pavez, and Gilles Louppe. Approximating likelihood ratios with calibrated
437 discriminative classifiers. *arXiv preprint arXiv:1506.02169*, 2015.
- 438 [36] Joeri Hermans, Volodimir Begy, and Gilles Louppe. Likelihood-free mcmc with approximate
439 likelihood ratios. *arXiv preprint arXiv:1903.04057*, 10, 2019.
- 440 [37] Conor Durkan, Iain Murray, and George Papamakarios. On contrastive learning for likelihood-
441 free inference. In *International conference on machine learning*, pages 2771–2781. PMLR,
442 2020.
- 443 [38] Owen Thomas, Ritabrata Dutta, Jukka Corander, Samuel Kaski, and Michael U. Gutmann.
444 Likelihood-free inference by ratio estimation, 2020. arXiv:1611.10242.
- 445 [39] Thomas Müller, Brian McWilliams, Fabrice Rousselle, Markus Gross, and Jan Novák. Neural
446 importance sampling. *ACM Transactions on Graphics (TOG)*, 38(5):1–19, 2019.
- 447 [40] Maximilian Dax, Stephen R. Green, Jonathan Gair, Michael Pürerer, Jonas Wildberger, Jakob H.
448 Macke, Alessandra Buonanno, and Bernhard Schölkopf. Neural Importance Sampling for
449 Rapid and Reliable Gravitational-Wave Inference. *Phys. Rev. Lett.*, 130(17):171403, 2023.
450 arXiv:2210.05686, doi:10.1103/PhysRevLett.130.171403.
- 451 [41] Michael S Albergo, Nicholas M Boffi, and Eric Vanden-Eijnden. Stochastic interpolants: A
452 unifying framework for flows and diffusions. *arXiv preprint arXiv:2303.08797*, 2023.
- 453 [42] Yang Song, Conor Durkan, Iain Murray, and Stefano Ermon. Maximum likelihood training of
454 score-based diffusion models. *Advances in Neural Information Processing Systems*, 34:1415–
455 1428, 2021.
- 456 [43] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks
457 for biomedical image segmentation. In *Medical Image Computing and Computer-Assisted
458 Intervention–MICCAI 2015: 18th International Conference, Munich, Germany, October 5-9,
459 2015, Proceedings, Part III 18*, pages 234–241. Springer, 2015.
- 460 [44] Yann N Dauphin, Angela Fan, Michael Auli, and David Grangier. Language modeling with
461 gated convolutional networks. In *International conference on machine learning*, pages 933–941.
462 PMLR, 2017.
- 463 [45] Jan-Matthis Lueckmann, Jan Boelts, David Greenberg, Pedro Goncalves, and Jakob Macke.
464 Benchmarking simulation-based inference. In *International Conference on Artificial Intelligence
465 and Statistics*, pages 343–351. PMLR, 2021.
- 466 [46] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
467 recognition, 2015. arXiv:1512.03385.
- 468 [47] Jerome H Friedman. On multivariate goodness-of-fit and two-sample testing. *Statistical
469 Problems in Particle Physics, Astrophysics, and Cosmology*, 1:311, 2003.
- 470 [48] David Lopez-Paz and Maxime Oquab. Revisiting classifier two-sample tests. *arXiv preprint
471 arXiv:1610.06545*, 2016.
- 472 [49] Arnaud Delaunoy, Joeri Hermans, François Rozet, Antoine Wehenkel, and Gilles Louppe.
473 Towards reliable simulation-based inference with balanced neural ratio estimation, 2022. arXiv:
474 2208.13624.
- 475 [50] Arnaud Delaunoy, Benjamin Kurt Miller, Patrick Forré, Christoph Weniger, and Gilles
476 Louppe. Balancing simulation-based inference for conservative posteriors. *arXiv preprint
477 arXiv:2304.10978*, 2023.
- 478 [51] R. Abbott et al. GWTC-3: Compact Binary Coalescences Observed by LIGO and Virgo
479 During the Second Part of the Third Observing Run. *arXiv preprint arXiv:2111.03606*, 11 2021.
480 arXiv:2111.03606.

- 481 [52] R. Abbott et al. Population Properties of Compact Objects from the Second LIGO-Virgo
482 Gravitational-Wave Transient Catalog. *Astrophys. J. Lett.*, 913(1):L7, 2021. arXiv:2010.
483 14533, doi:10.3847/2041-8213/abe949.
- 484 [53] B. P. Abbott et al. GW170817: Measurements of neutron star radii and equation of state. *Phys.*
485 *Rev. Lett.*, 121(16):161101, 2018. arXiv:1805.11581, doi:10.1103/PhysRevLett.121.
486 161101.
- 487 [54] R. Abbott et al. Tests of general relativity with binary black holes from the second LIGO-Virgo
488 gravitational-wave transient catalog. *Phys. Rev. D*, 103(12):122002, 2021. arXiv:2010.14529,
489 doi:10.1103/PhysRevD.103.122002.
- 490 [55] B. P. Abbott et al. A gravitational-wave standard siren measurement of the Hubble constant.
491 *Nature*, 551(7678):85–88, 2017. arXiv:1710.05835, doi:10.1038/nature24471.
- 492 [56] J. Veitch, V. Raymond, B. Farr, W Farr, P. Graff, S. Vitale, et al. Parameter estimation for
493 compact binaries with ground-based gravitational-wave observations using the LALInference
494 software library. *Phys. Rev.*, D91(4):042003, 2015. arXiv:1409.7215, doi:10.1103/
495 PhysRevD.91.042003.
- 496 [57] Gregory Ashton et al. BILBY: A user-friendly Bayesian inference library for gravitational-
497 wave astronomy. *Astrophys. J. Suppl.*, 241(2):27, 2019. arXiv:1811.02042, doi:10.3847/
498 1538-4365/ab06fc.
- 499 [58] I. M. Romero-Shaw et al. Bayesian inference for compact binary coalescences with bilby:
500 validation and application to the first LIGO–Virgo gravitational-wave transient catalogue. *Mon.*
501 *Not. Roy. Astron. Soc.*, 499(3):3295–3319, 2020. arXiv:2006.00714, doi:10.1093/mnras/
502 staa2850.
- 503 [59] Joshua S Speagle. dynesty: a dynamic nested sampling package for estimating Bayesian poste-
504 riors and evidences. *Monthly Notices of the Royal Astronomical Society*, 493(3):3132–3158,
505 Feb 2020. URL: <http://dx.doi.org/10.1093/mnras/staa278>, arXiv:1904.02180,
506 doi:10.1093/mnras/staa278.
- 507 [60] Nicholas Metropolis, Arianna W Rosenbluth, Marshall N Rosenbluth, Augusta H Teller, and
508 Edward Teller. Equation of state calculations by fast computing machines. *The journal of*
509 *chemical physics*, 21(6):1087–1092, 1953.
- 510 [61] W. K. Hastings. Monte Carlo sampling methods using Markov chains and their applica-
511 tions. *Biometrika*, 57(1):97–109, 04 1970. arXiv:[https://academic.oup.com/biomet/
512 article-pdf/57/1/97/23940249/57-1-97.pdf](https://academic.oup.com/biomet/article-pdf/57/1/97/23940249/57-1-97.pdf), doi:10.1093/biomet/57.1.97.
- 513 [62] John Skilling. Nested sampling for general Bayesian computation. *Bayesian Analysis*, 1(4):833
514 – 859, 2006. doi:10.1214/06-BA127.
- 515 [63] Elena Cuoco, Jade Powell, Marco Cavaglià, Kendall Ackley, Michał Bejger, Chayan Chatterjee,
516 Michael Coughlin, Scott Coughlin, Paul Easter, Reed Essick, et al. Enhancing gravitational-
517 wave science with machine learning. *Machine Learning: Science and Technology*, 2(1):011002,
518 5 2020. arXiv:2005.03745, doi:10.1088/2632-2153/abb93a.
- 519 [64] Hunter Gabbard, Chris Messenger, Ik Siong Heng, Francesco Tonolini, and Roderick
520 Murray-Smith. Bayesian parameter estimation using conditional variational autoencoders
521 for gravitational-wave astronomy. *Nature Phys.*, 18(1):112–117, 2022. arXiv:1909.06296,
522 doi:10.1038/s41567-021-01425-7.
- 523 [65] Stephen R. Green, Christine Simpson, and Jonathan Gair. Gravitational-wave parameter
524 estimation with autoregressive neural network flows. *Phys. Rev. D*, 102(10):104057, 2020.
525 arXiv:2002.07656, doi:10.1103/PhysRevD.102.104057.
- 526 [66] Arnaud Delaunoy, Antoine Wehenkel, Tanja Hinderer, Samaya Nissanke, Christoph Weniger,
527 Andrew R. Williamson, and Gilles Louppe. Lightning-Fast Gravitational Wave Parameter
528 Inference through Neural Amortization. In *Third Workshop on Machine Learning and the*
529 *Physical Sciences*, 10 2020. arXiv:2010.12931.

- 530 [67] Stephen R. Green and Jonathan Gair. Complete parameter inference for GW150914 using deep
531 learning. *Mach. Learn. Sci. Tech.*, 2(3):03LT01, 2021. arXiv:2008.03312, doi:10.1088/
532 2632-2153/abfaed.
- 533 [68] Maximilian Dax, Stephen R. Green, Jonathan Gair, Michael Deistler, Bernhard Schölkopf, and
534 Jakob H. Macke. Group equivariant neural posterior estimation. In *International Conference on*
535 *Learning Representations*, 11 2022. arXiv:2111.13139.
- 536 [69] Chayan Chatterjee, Linqing Wen, Damon Beveridge, Foivos Diakogiannis, and Kevin Vinsen.
537 Rapid localization of gravitational wave sources from compact binary coalescences using deep
538 learning. *arXiv preprint arXiv:2207.14522*, 7 2022. arXiv:2207.14522.
- 539 [70] Stefan T. Radev, Ulf K. Mertens, Andreass Voss, Lynton Ardizzone, and Ullrich Köthe.
540 Bayesflow: Learning complex stochastic models with invertible neural networks, 2020.
541 arXiv:2003.06281.
- 542 [71] Conor Durkan, Artur Bekasov, Iain Murray, and George Papamakarios. Neural spline flows.
543 In *Advances in Neural Information Processing Systems*, pages 7509–7520, 2019. arXiv:
544 1906.04032.
- 545 [72] B.P. Abbott et al. Observation of Gravitational Waves from a Binary Black Hole Merger.
546 *Phys. Rev. Lett.*, 116(6):061102, 2016. arXiv:1602.03837, doi:10.1103/PhysRevLett.
547 116.061102.
- 548 [73] Jonas Wildberger, Maximilian Dax, Stephen R. Green, Jonathan Gair, Michael Pürrer, Jakob H.
549 Macke, Alessandra Buonanno, and Bernhard Schölkopf. Adapting to noise distribution shifts
550 in flow-based gravitational-wave inference. *Phys. Rev. D*, 107(8):084046, 2023. arXiv:
551 2211.08801, doi:10.1103/PhysRevD.107.084046.
- 552 [74] Taco Cohen and Max Welling. Group equivariant convolutional networks. In Maria-Florina
553 Balcan and Kilian Q. Weinberger, editors, *Proceedings of the 33rd International Conference*
554 *on Machine Learning, ICML 2016, New York City, NY, USA, June 19-24, 2016*, volume 48
555 of *JMLR Workshop and Conference Proceedings*, pages 2990–2999. JMLR.org, 2016. URL:
556 <http://proceedings.mlr.press/v48/cohen16.html>.
- 557 [75] Mark Hannam, Patricia Schmidt, Alejandro Bohé, Leïla Haegel, Sascha Husa, Frank Ohme,
558 Geraint Pratten, and Michael Pürrer. Simple model of complete precessing black-hole-
559 binary gravitational waveforms. *Phys. Rev. Lett.*, 113:151101, Oct 2014. URL: [https://](https://link.aps.org/doi/10.1103/PhysRevLett.113.151101)
560 link.aps.org/doi/10.1103/PhysRevLett.113.151101, doi:10.1103/PhysRevLett.
561 113.151101.
- 562 [76] Sebastian Khan, Sascha Husa, Mark Hannam, Frank Ohme, Michael Pürrer, Xisco Jiménez
563 Forteza, and Alejandro Bohé. Frequency-domain gravitational waves from nonprecessing
564 black-hole binaries. II. A phenomenological model for the advanced detector era. *Phys. Rev.*,
565 D93(4):044007, 2016. arXiv:1508.07253, doi:10.1103/PhysRevD.93.044007.
- 566 [77] Alejandro Bohé, Mark Hannam, Sascha Husa, Frank Ohme, Michael Pürrer, and Patricia
567 Schmidt. PhenomPv2 – technical notes for the LAL implementation. *LIGO Technical Document*,
568 *LIGO-T1500602-v4*, 2016. URL: <https://dcc.ligo.org/LIGO-T1500602/public>.
- 569 [78] Benjamin Farr, Evan Ochsner, Will M. Farr, and Richard O’Shaughnessy. A more effective
570 coordinate system for parameter estimation of precessing compact binaries from gravitational
571 waves. *Phys. Rev. D*, 90(2):024018, 2014. arXiv:1404.7070, doi:10.1103/PhysRevD.90.
572 024018.
- 573 [79] J.R. Dormand and P.J. Prince. A family of embedded runge-kutta formulae. *Journal of*
574 *Computational and Applied Mathematics*, 6(1):19–26, 1980. URL: [https://](https://www.sciencedirect.com/science/article/pii/0771050X80900133)
575 www.sciencedirect.com/science/article/pii/0771050X80900133, doi:[https://](https://doi.org/10.1016/0771-050X(80)90013-3)
576 [doi.org/10.1016/0771-050X\(80\)90013-3](https://doi.org/10.1016/0771-050X(80)90013-3).

577 **A Gaussian flow**

578 We here derive the form of a vector field $v_t(\theta)$ that restricts the resulting continuous flow to a one
 579 dimensional Gaussian with mean $\hat{\mu}$ variance $\hat{\sigma}^2$. With the optimal transport path $\mu_t(\theta) = t\theta_1$,
 580 $\sigma_t(\theta) = 1 - (1 - \sigma_{\min})t \equiv \sigma_t$ from [16], the sample-conditional probability path (5) reads

$$p_t(\theta|\theta_1) = \mathcal{N}[t\theta_1, \sigma_t^2](\theta). \quad (10)$$

581 We set our target distribution

$$q_1(\theta_1) = \mathcal{N}[\hat{\mu}, \hat{\sigma}^2](\theta_1). \quad (11)$$

582 To derive the marginal probability path and the marginal vector field we need two identities for
 583 the convolution $*$ of Gaussian densities. Recall that the convolution of two function is defined by
 584 $f * g(x) = \int f(x - y)g(y) dy$. We define the function

$$g_{\mu, \sigma^2}(\theta) = \theta \cdot \mathcal{N}[\mu, \sigma^2](\theta). \quad (12)$$

585 Then the following holds

$$\mathcal{N}[\mu_1, \sigma_1^2] * \mathcal{N}[\mu_2, \sigma_2^2] = \mathcal{N}[\mu_1 + \mu_2, \sigma_1^2 + \sigma_2^2] \quad (13)$$

$$g_{0, \sigma_1^2} * \mathcal{N}[\mu_2, \sigma_2^2] = \frac{\sigma_1^2}{\sigma_1^2 + \sigma_2^2} \left(g_{\mu_2, \sigma_1^2 + \sigma_2^2} - \mu_2 \mathcal{N}[\mu_2, \sigma_1^2 + \sigma_2^2] \right) \quad (14)$$

586 **Marginal probability paths**

587 Marginalizing over θ_1 in (10) with (11), we find

$$\begin{aligned} p_t(\theta) &= \int p_t(\theta|\theta_1)q(\theta_1) d\theta_1 \\ &= \int \mathcal{N}[t\theta_1, \sigma_t^2](\theta) \mathcal{N}[\hat{\mu}, \hat{\sigma}^2](\theta_1) d\theta_1 \\ &= \int \mathcal{N}[0, \sigma_t^2](\theta - t\theta_1) \mathcal{N}(t\hat{\mu}, (t\hat{\sigma})^2)(t\theta_1) \cdot t d\theta_1 \\ &= \int \mathcal{N}[0, \sigma_t^2](\theta - \theta_1^t) \mathcal{N}[t\hat{\mu}, (t\hat{\sigma})^2](\theta_1^t) d\theta_1^t \\ &= \mathcal{N}[t\hat{\mu}, \sigma_t^2 + (t\hat{\sigma})^2](\theta) \end{aligned} \quad (15)$$

588 where we defined $\theta_1^t = t\theta_1$ and used (13).

589 **Marginal vector field**

590 We now calculate the marginalized vector field $u_t(\theta)$ based on equation (8) in [16]. Using the
 591 sample-conditional vector field (6) and the distributions (10) and (11) we find

$$\begin{aligned} u_t(\theta) &= \int u_t(\theta|\theta_1) \frac{p_t(\theta|\theta_1)q(\theta_1)}{p_t(\theta)} d\theta_1 \\ &= \frac{1}{p_t(\theta)} \int \frac{(\theta_1 - (1 - \sigma_{\min})\theta)}{\sigma_t} \cdot \mathcal{N}[t\theta_1, \sigma_t^2](\theta) \cdot \mathcal{N}[\hat{\mu}, \hat{\sigma}^2](\theta_1) d\theta_1 \\ &= \frac{1}{p_t(\theta)} \int \frac{(\theta_1 - (1 - \sigma_{\min})\theta)}{\sigma_t} \cdot \mathcal{N}[0, \sigma_t^2](\theta - t\theta_1) \cdot \mathcal{N}[t\hat{\mu}, (t\hat{\sigma})^2](t\theta_1) \cdot t d\theta_1 \\ &= \frac{1}{p_t(\theta)} \int \frac{(\theta_1^t - (1 - \sigma_{\min})t \cdot \theta)}{\sigma_t \cdot t} \cdot \mathcal{N}[0, \sigma_t^2](\theta - \theta_1^t) \cdot \mathcal{N}[t\hat{\mu}, (t\hat{\sigma})^2](\theta_1^t) \cdot d\theta_1^t \\ &= \frac{1}{p_t(\theta)} \int \frac{(-\theta_1'' + (1 - (1 - \sigma_{\min})t) \cdot \theta)}{\sigma_t \cdot t} \cdot \mathcal{N}[0, \sigma_t^2](\theta_1'') \cdot \mathcal{N}[t\hat{\mu}, (t\hat{\sigma})^2](\theta - \theta_1'') \cdot d\theta_1'' \\ &= \frac{1}{p_t(\theta)} \int \frac{(-\theta_1'' + \sigma_t \cdot \theta)}{\sigma_t \cdot t} \cdot \mathcal{N}[0, \sigma_t^2](\theta_1'') \cdot \mathcal{N}[t\hat{\mu}, (t\hat{\sigma})^2](\theta - \theta_1'') \cdot d\theta_1'' \end{aligned} \quad (16)$$

592 where we used the change of variables $\theta'_1 = t\theta_1$ and $\theta''_1 = \theta - \theta'_1$. Now we evaluate this expression
 593 using (12), then the identities (13) and (14) and the marginal probability (15)

$$\begin{aligned}
 u_t(\theta) &= \frac{-1}{p_t(\theta) \cdot \sigma_t \cdot t} \left(g_{0, \sigma_t^2} * \mathcal{N}[t\hat{\mu}, (t\hat{\sigma})^2] \right) (\theta) + \frac{\theta}{p_t(\theta) \cdot t} \left(\mathcal{N}[0, \sigma_t^2] * \mathcal{N}[t\hat{\mu}, (t\hat{\sigma})^2] \right) (\theta) \\
 &= \frac{-1}{p_t(\theta) \cdot \sigma_t \cdot t} \frac{(\theta - t\hat{\mu}) \cdot \sigma_t^2}{\sigma_t^2 + (t\hat{\sigma})^2} \cdot \mathcal{N}[t\hat{\mu}, (\sigma_t^2 + (t\hat{\sigma})^2)] (\theta) + \frac{\theta}{p_t(\theta) \cdot t} \mathcal{N}[t\hat{\mu}, (\sigma_t^2 + (t\hat{\sigma})^2)] (\theta) \\
 &= \frac{(\sigma_t^2 + (t\hat{\sigma})^2)\theta - (\theta - t\hat{\mu}) \cdot \sigma_t}{p_t(\theta) \cdot t \cdot (\sigma_t^2 + (t\hat{\sigma})^2)} \cdot p_t(\theta) \\
 &= \frac{(\sigma_t^2 + (t\hat{\sigma})^2) - \sigma_t}{t \cdot (\sigma_t^2 + (t\hat{\sigma})^2)} \cdot \theta + \frac{t\hat{\mu} \cdot \sigma_t}{t \cdot (\sigma_t^2 + (t\hat{\sigma})^2)}.
 \end{aligned} \tag{17}$$

594 By choosing a vector field v_t of the form (17) with learnable parameters $\hat{\mu}, \hat{\sigma}^2$, we can thus define a
 595 continuous flow that is restricted to a one dimensional Gaussian.

596 B Mass covering properties of flows

597 In this supplement, we investigate the mass covering properties of continuous normalizing flows
 598 trained using mean squared error and in particular prove Theorem 1. We first recall the notation from
 599 the main part. We always assume that the data is distributed according to $p_1(\theta)$. In addition, there is a
 600 known and simple base distribution p_0 and we assume that there is a vector field $u_t : [0, 1] \times \mathbb{R}^d \rightarrow \mathbb{R}^d$
 601 that connects p_0 and p_1 in the following sense. We denote by ϕ_t the flow generated by u_t , i.e., ϕ_t
 602 satisfies

$$\partial_t \phi_t(\theta) = u_t(\phi_t(\theta)). \tag{18}$$

603 Then we assume that $(\phi_1)_* p_0 = p_1$ and we also define the interpolations $p_t = (\phi_t)_* p_0$.

604 We do not have access to the ground truth distributions p_t and the vector field u_t but we try to learn a
 605 vector field v_t approximating u_t . We denote its flow by ψ_t and we define $q_t = (\psi_t)_* q_0$ and $q_0 = p_0$.
 606 We are interested in the mass covering properties of the learned approximation q_1 of p_1 . In particular,
 607 we want to relate the KL-divergence $D_{\text{KL}}(p_1 || q_1)$ to the mean squared error,

$$\text{MSE}_p(u, v) = \int_0^1 dt \int p_t(d\theta) (u_t(\theta) - v_t(\theta))^2, \tag{19}$$

608 of the generating vector fields. The first observation is that without any regularity assumptions on v_t
 609 it is impossible to obtain any bound on the KL-divergence in terms of the mean squared error.

610 **Lemma 1.** *For every $\varepsilon > 0$ there are vector field u_t and v_t and a base distribution $p_0 = q_0$ such that*

$$\text{MSE}_p(u, v) < \varepsilon \text{ and } D_{\text{KL}}(p_1 || q_1) = \infty. \tag{20}$$

611 *In addition we can construct u_t and v_t such that the support of p_1 is larger than the support of q_1 .*

612 *Proof.* We consider the uniform distribution $p_0 = q_0 \sim \mathcal{U}([-1, 1])$ and the vector fields

$$u_t(\theta) = 0 \tag{21}$$

613 and

$$v_t(\theta) = \begin{cases} \varepsilon & \text{for } 0 \leq \theta < \varepsilon, \\ 0 & \text{otherwise.} \end{cases} \tag{22}$$

614 As before, let ϕ_t denote the flow of the vector field u_t and similarly ψ_t denote the flow of v_t . Clearly
 615 $\phi_t(\theta) = \theta$. On the other hand

$$\psi_t(\theta) = \begin{cases} \min(\theta + \varepsilon t, \varepsilon) & \text{if } 0 \leq \theta < \varepsilon, \\ \theta & \text{otherwise.} \end{cases} \tag{23}$$

616 In particular

$$\psi_1(\theta) = \begin{cases} \varepsilon & \text{if } 0 \leq \theta < \varepsilon, \\ \theta & \text{otherwise.} \end{cases} \quad (24)$$

617 This implies that $p_1 = (\phi_1)_* p_0 \sim \mathcal{U}([-1, 1])$. On the other hand $q_1 = (\psi_1)_* q_0$ has support in
 618 $[-1, 0] \cup [\varepsilon, 1]$. In particular, the distribution of q_1 is not mass covering with respect to p_1 and
 619 $D_{\text{KL}}(p_1 || q_1) = \infty$. Finally, we observe that the MSE can be arbitrarily small

$$\text{MSE}_p(u, v) = \int_0^1 dt \int p_t(d\theta) |u_t(\theta) - v_t(\theta)|^2 = \int_0^1 \int_0^\varepsilon \frac{1}{2} \varepsilon^2 = \frac{\varepsilon^3}{2}. \quad (25)$$

620 Here we used that the density of $p_t(d\theta)$ is $1/2$ for $-1 \leq \theta \leq 1$. \square

621 We see that an arbitrary small MSE-loss cannot ensure that the probability distribution is mass
 622 covering and the KL-divergence is finite. On a high level this can be explained by the fact that
 623 for vector fields v_t that are not Lipschitz continuous the flow is not necessarily continuous, and
 624 we can generate holes in the distribution. Note that we chose p_0 to be a uniform distribution for
 625 simplicity, but the result extends to any smooth distribution, in particular the result does not rely on
 626 the discontinuity of p_0 .

627 Next, we investigate the mass covering property for Lipschitz continuous flows. When the flows u_t
 628 and v_t are Lipschitz continuous (in θ) this ensures that the flows ψ_1 and ϕ_1 are continuous in x and it
 629 is not possible to create holes in the distribution as shown above for non-continuous vector fields. We
 630 show a weaker bound in this setting.

631 **Lemma 2.** *For every $0 \leq \delta \leq 1$ there is a base distribution $p_0 = q_0$ and the are Lipschitz-continuous
 632 vector fields u_t and v_t such that $\text{MSE}_p(u, v) = \delta$ and*

$$D_{\text{KL}}(p_1 || q_1) \geq \frac{1}{2} \text{MSE}_p(u, v)^{1/3}. \quad (26)$$

633 *Proof.* We consider p_0, q_0 and u_t as in Lemma 1, and we define

$$v_t(\theta) = \begin{cases} 2\theta & \text{for } 0 \leq \theta < \varepsilon, \\ 2\varepsilon - \theta & \text{for } \varepsilon \leq \theta < 2\varepsilon, \\ 0 & \text{otherwise.} \end{cases} \quad (27)$$

634 Then we can calculate for $0 \leq \theta \leq e^{-2}\varepsilon$ that

$$\psi_t(\theta) = \theta e^{2t}. \quad (28)$$

635 Similarly we obtain for $\varepsilon \leq \theta \leq 2\varepsilon$ (solving the ODE $f' = 2f$)

$$\psi_t(\theta) = 2\varepsilon - (\theta - \varepsilon)e^{-2t}. \quad (29)$$

636 We find

$$\psi_1(0) = 0, \psi_1(e^{-2}\varepsilon) = \varepsilon, \psi_1(\varepsilon) = 2 - \varepsilon e^{-2}; \psi_2(2\varepsilon) = 2\varepsilon. \quad (30)$$

637 Next we find for the densities of q_1 that

$$q_1(\psi_1(\theta)) = q_0(\theta) |\psi_1'(\theta)|^{-1} = \frac{1}{2} \begin{cases} e^{-2} & \text{for } 0 \leq \theta \leq e^{-2}\varepsilon, \\ e^2 & \text{for } \varepsilon \leq \theta \leq 2\varepsilon. \end{cases} \quad (31)$$

638 Together with (30) this implies that the density of q_1 is given by

$$q_1(\theta) = \frac{1}{2} \begin{cases} e^{-2} & \text{for } 0 \leq \theta \leq \varepsilon, \\ e^2 & \text{for } 2\varepsilon - \varepsilon e^{-2} \leq \theta \leq 2\varepsilon. \end{cases} \quad (32)$$

639 Note that $p_1(\theta) = 1/2$ for $-1 \leq \theta \leq 1$ and therefore

$$\int_0^\varepsilon \ln \frac{p_1(\theta)}{q_1(\theta)} p_1(d\theta) = \int_0^\varepsilon \ln(e^2) \frac{1}{2} d\theta = \varepsilon, \quad (33)$$

640 and

$$\int_{2\varepsilon - \varepsilon e^{-2}}^{2\varepsilon} \ln \frac{p_1(\theta)}{q_1(\theta)} p_1(d\theta) = \int_{2\varepsilon - \varepsilon e^{-2}}^{2\varepsilon} \ln(e^{-2}) \frac{1}{2} d\theta = -\varepsilon e^{-2}. \quad (34)$$

641 Moreover we note

$$\int_{\varepsilon}^{2\varepsilon - \varepsilon e^{-2}} q_1(d\varepsilon) = \int_{e^{-2\varepsilon}}^{\varepsilon} q_0(d\varepsilon) = \frac{1}{2} \varepsilon (1 - e^{-2}) = \int_{\varepsilon}^{2\varepsilon - \varepsilon e^{-2}} p_1(d\varepsilon), \quad (35)$$

642 which implies (by positivity of the KL-divergence) that

$$\int_{\varepsilon}^{2\varepsilon - \varepsilon e^{-2}} \ln \left(\frac{p_1(\theta)}{q_1(\theta)} \right) p_1(d\theta) \geq 0. \quad (36)$$

643 We infer using also $p_1(\theta) = q_1(\theta) = 1/2$ for $\theta \in [-1, 0] \cap [2\varepsilon, 1]$ that

$$D_{\text{KL}}(p_1 || q_1) = \int \ln \left(\frac{p_1(\theta)}{q_1(\theta)} \right) p_1(d\theta) \geq \varepsilon (1 - e^{-2}). \quad (37)$$

644 On the other hand we can bound

$$\int_0^1 dt \int p_t(d\theta) |v_t(\theta) - u_t(\theta)|^2 = \frac{1}{2} \int_0^1 dt \int_0^{2\varepsilon} |u_t(\theta)|^2 = \int_0^{\varepsilon} s^2 ds = \frac{\varepsilon^3}{3}. \quad (38)$$

645 We conclude that

$$D_{\text{KL}}(p_1 || q_1) \geq \frac{1}{2} (\text{MSE}_p(u, v))^{1/3}. \quad (39)$$

646 In particular, it is not possible to bound the KL-divergence by the MSE even when the vector fields
647 are Lipschitz continuous. \square

648 Let us put this into context. It was already shown in [41] that we can, in general, not bound the
649 forward KL-divergence by the mean squared error and our Lemmas 1 and 2 are concrete examples.
650 On the other hand, when considering SDEs the KL-divergence can be bounded by the mean squared
651 error of the drift terms as shown in [42]. Indeed, in [41] the favorable smoothing effect was carefully
652 investigated.

653 Here we show that we can alternatively obtain an upper bound on the KL-divergence when assuming
654 that u_t , v_t , and p_0 satisfy additional regularity assumptions. This allows us to recover the mass
655 covering property from bounds on the means squared error for sufficiently smooth vector fields. The
656 scaling is nevertheless still weaker than for SDEs.

657 We now state our assumptions. We denote the gradient with respect to θ by $\nabla = \nabla_{\mu}$ and second
658 derivatives by $\nabla^2 = \nabla_{\mu\nu}^2$. When applying the chain rule, we leave the indices implicit. We denote by
659 $|\cdot|$ the Frobenius norm $|A| = \left(\sum_{ij} A_{ij}^2 \right)^{1/2}$ of a matrix. The Frobenius norm is submultiplicative,
660 i.e., $|AB| \leq |A| \cdot |B|$ and directly generalizes to higher order tensors.

661 *Assumption 1.* We assume that

$$|\nabla u_t| \leq L, \quad |\nabla v_t| \leq L, \quad |\nabla^2 u_t| \leq L', \quad |\nabla^2 v_t| \leq L'. \quad (40)$$

662 We require one further assumption on p_0 .

663 *Assumption 2.* There is a constant C_1 such that

$$|\nabla \ln p_0(\theta)| \leq C_1 (1 + |\theta|). \quad (41)$$

664 We also assume that

$$\mathbb{E}_{p_0} |\theta|^2 < C_2 < \infty. \quad (42)$$

665 Note that (41) holds, e.g., if p_0 follows a Gaussian distribution but also for smooth distribution with
666 slower decay at ∞ . If we assume that $|\nabla \ln p_0(\theta)|$ is bounded the proof below simplifies slightly.
667 This is, e.g., the case if $p_0(\theta) \sim e^{-|\theta|}$ as $|\theta| \rightarrow \infty$.

668 We need some additional notation. It is convenient to introduce $\phi_t^s = \phi_t \circ (\phi_s)^{-1}$, i.e., the flow from
669 time s to t (in particular $\phi_t^0 = \phi_t$) and similarly for ψ . We can now restate and prove Theorem 1.

670 **Theorem 2.** Let $p_0 = q_0$ and assume u_t and v_t are two vector fields whose flows satisfy $p_1 = (\phi_1)_* p_0$
671 and $q_1 = (\psi_1)_* q_0$. Assume that p_0 satisfies Assumption 2 and u_t and v_t satisfy Assumption 1. Then
672 there is a constant $C > 0$ depending on L, L', C_1, C_2 , and d such that (for $\text{MSE}_p(u, v) < 1$)

$$D_{\text{KL}}(p_1 || q_1) \leq C \text{MSE}_p(u, v)^{\frac{1}{2}}. \quad (43)$$

673 *Remark 1.* We do not claim that our results are optimal, it might be possible to find similar bounds
674 for the forward KL-divergence with weaker assumptions. However, we emphasize that Lemma 2
675 shows that the result of the theorem is not true without the assumption on the second derivative of v_t
676 and u_t .

677 *Proof.* We want to control $D_{\text{KL}}(p_1 || q_1)$. It can be shown that (see equation above (25) in [42] or
678 Lemma 2.19 in [41])

$$\partial_t D_{\text{KL}}(p_t || q_t) = - \int p_t(d\theta) (u_t(\theta) - v_t(\theta)) \cdot (\nabla \ln p_t(\theta) - \nabla \ln q_t(\theta)). \quad (44)$$

679 Using Cauchy-Schwarz we can bound this by

$$\partial_t D_{\text{KL}}(p_t || q_t) \leq \left(\int p_t(d\theta) |u_t(\theta) - v_t(\theta)|^2 \right)^{\frac{1}{2}} \left(\int p_t(d\theta) |\nabla \ln p_t(\theta) - \nabla \ln q_t(\theta)|^2 \right)^{\frac{1}{2}}. \quad (45)$$

680 We use the relation (see (3))

$$\ln(p_t(\phi_t(\theta_0))) = \ln(p_0(\theta_0)) - \int_0^t (\text{div } u_s)(\phi_s(\theta_0)) ds, \quad (46)$$

681 which can be equivalently rewritten (setting $\theta = \phi_t \theta_0$) as

$$\ln(p_t(\theta)) = \ln(p_0(\phi_0^t \theta)) - \int_0^t (\text{div } u_s)(\phi_s^t \theta) ds. \quad (47)$$

682 We use the following relation for $\nabla \phi_s^t$

$$\nabla \phi_s^t(\theta) = \exp \left(\int_t^s d\tau (\nabla u_\tau)(\phi_\tau^t(\theta)) \right). \quad (48)$$

683 This relation is standard and can be directly deduced from the following ODE for $\nabla \phi_s^t$

$$\partial_s \nabla \phi_s^t(\theta) = \nabla \partial_s \phi_s^t(\theta) = \nabla (u_s(\phi_s^t(\theta))) = ((\nabla u_s)(\phi_s^t(\theta))) \cdot \nabla \phi_s^t(\theta). \quad (49)$$

684 We can conclude that for $0 \leq s, t \leq 1$ the bound

$$|\nabla \phi_s^t(\theta)| \leq e^L \quad (50)$$

685 holds. We find

$$\begin{aligned} |\nabla \ln(p_t(\theta))| &= \left| \nabla \ln(p_0)(\phi_0^t \theta) \cdot \nabla \phi_0^t(\theta) - \int_0^t (\nabla \text{div } u_s)(\phi_s^t \theta) \cdot \nabla \phi_s^t(\theta) ds \right| \\ &\leq |\nabla \ln(p_0)(\phi_0^t \theta)| e^L + L' e^L, \end{aligned} \quad (51)$$

686 and a similar bound holds for q_t . In words, we have shown that the score of p_t at θ can be bounded by
687 the score of p_0 of theta transported along the vector field u_t minus a correction which quantifies the
688 change of score along the path. We now bound using the definition $p_t = (\phi_t)_* p_0$ and the assumption
689 (41)

$$\begin{aligned} \int p_t(d\theta) |\nabla \ln p_0(\phi_0^t(\theta))|^2 &= \int p_0(d\theta_0) |\nabla \ln p_0(\phi_0^t \phi_t(\theta_0))|^2 = \mathbb{E}_{p_0} |\nabla \ln p_0(\theta_0)|^2 \\ &\leq \mathbb{E}_{p_0} (C_1(1 + |\theta_0|)^2) \leq 2C_1^2(1 + \mathbb{E}_{p_0} |\theta_0|^2) \leq 2C_1^2(1 + C_2^2). \end{aligned} \quad (52)$$

690 Similarly we obtain using $q_0 = p_0$

$$\int p_t(d\theta) |\nabla \ln q_0(\psi_0^t \theta)|^2 = \int p_0(d\theta_0) |\nabla \ln q_0(\psi_0^t \phi_t \theta_0)|^2. \quad (53)$$

691 In words, to control the score of q integrated with respect to p_t we need to control the distortion we
 692 obtain when moving forward with u and backwards with v . We investigate $\psi_0^t \phi_t(\theta_0)$. We find

$$\partial_h \psi_t^{t+h} \phi_{t+h}^t(\theta)|_{h=0} = u_t(\theta) - v_t(\theta). \quad (54)$$

693 This implies

$$\partial_t(\psi_0^t \phi_t)(\theta_0) = \partial_h(\psi_0^t \psi_t^{t+h} \phi_{t+h}^t)(\theta_0)|_{h=0} = (\nabla \psi_0^t)(\phi_t(\theta_0)) \cdot ((u_t - v_t)(\phi_t(\theta_0))). \quad (55)$$

694 Using (50) we conclude that

$$\begin{aligned} |\psi_0^t \phi_t(\theta_0) - \theta_0| &\leq \left| \int_0^t \partial_s \psi_0^s \phi_s(\theta_0) \, ds \right| \leq \int_0^t |(\nabla \psi_0^s)(\phi_s(\theta_0))| \cdot |u_s - v_s|(\phi_s(\theta_0)) \, ds \\ &\leq e^L \int_0^t |u_s - v_s|(\phi_s(\theta_0)) \, ds. \end{aligned} \quad (56)$$

695 We use this and the assumption (41) to continue to estimate (53) as follows

$$\begin{aligned} \int p_t(d\theta) |\nabla \ln q_0(\psi_0^t \theta)|^2 &= \int p_0(d\theta_0) |\nabla \ln q_0(\psi_0^t \phi_t(\theta_0))|^2 \\ &\leq C_1^2 \int p_0(d\theta_0) (1 + |\psi_0^t \phi_t(\theta_0)|)^2 \\ &\leq C_1^2 \int p_0(d\theta_0) (1 + |\psi_0^t \phi_t(\theta_0) - \theta_0| + |\theta_0|)^2 \\ &\leq 3C_1^2 + 3C_1^2 \int p_0(d\theta_0) (|\psi_0^t \phi_t(\theta_0) - \theta_0|^2 + |\theta_0|^2) \\ &\leq 3C_1^2 (1 + \mathbb{E}_{p_0} |\theta_0|^2) + 3C_1^2 e^{2L} \int p_0(d\theta_0) \left(\int_0^t ds |u_s - v_s|(\phi_s(\theta_0)) \right)^2. \end{aligned} \quad (57)$$

696 Here we used $(a + b + c)^2 \leq 3(a^2 + b^2 + c^2)$ in the second to last step. We bound the remaining
 697 integral using Cauchy-Schwarz as follows

$$\begin{aligned} \int p_0(d\theta_0) \left(\int_0^t |u_s - v_s|(\phi_s(\theta_0)) \, ds \right)^2 &\leq \int p_0(d\theta_0) \left(\int_0^t ds |u_s - v_s|^2(\phi_s(\theta_0)) \right) \left(\int_0^t ds 1^2 \right) \\ &\leq t \int_0^t ds \int p_0(d\theta_0) |u_s - v_s|^2(\phi_s(\theta_0)) \\ &= t \int_0^t ds \int p_s(d\theta_s) |u_s - v_s|^2(\theta_s) \\ &\leq \int_0^1 ds \int p_s(d\theta_s) |u_s - v_s|^2(\theta_s) = \text{MSE}_p(u, v). \end{aligned} \quad (58)$$

698 The last displays together imply

$$\int p_t(d\theta) |\nabla \ln q_0(\psi_0^t \theta)|^2 \leq 3C_1^2 (1 + \mathbb{E}_{p_0} |\theta_0|^2 + e^{2L} \text{MSE}_p(u, v)). \quad (59)$$

699 Now we have all the necessary ingredients to bound the derivative of the KL-divergence. We control
 700 the second integral in (45) using (51) (and again $(\sum_{i=1}^4 a_i)^2 \leq 4 \sum a_i^2$) as follows,

$$\begin{aligned} \int p_t(d\theta) |\nabla \ln p_t(\theta) - \nabla \ln q_t(\theta)|^2 \\ \leq 2 \cdot 2^2 \cdot L'^2 e^{2L} + 4e^{2L} \int p_t(d\theta) (|\nabla \ln q_0(\psi_0^t \theta)|^2 + |\nabla \ln p_0(\phi_0^t \theta)|^2). \end{aligned} \quad (60)$$

701 Using (52) and (59) we finally obtain

$$\begin{aligned} \int p_t(d\theta) |\nabla \ln p_t(\theta) - \nabla \ln q_t(\theta)|^2 &\leq 8 \cdot L'^2 e^{2L} + C_1^2 e^{2L} (20(1 + C_2^2) + 12 \text{MSE}_p(u, v)) \\ &\leq C(1 + \text{MSE}_p(u, v)) \end{aligned} \quad (61)$$

702 for some constant $C > 0$. Finally, we obtain

$$\begin{aligned}
 D_{\text{KL}}(p_1||q_1) &= \int_0^1 dt \partial_t D_{\text{KL}}(p_t||q_t) \\
 &\leq (C(1 + \text{MSE}_p(u, v)))^{\frac{1}{2}} \int_0^1 dt \left(\int p_t(d\theta) |u_t(\theta) - v_t(\theta)|^2 \right)^{\frac{1}{2}} \\
 &\leq (C(1 + \text{MSE}_p(u, v)))^{\frac{1}{2}} \left(\int_0^1 dt \int p_t(d\theta) |u_t(\theta) - v_t(\theta)|^2 \right)^{\frac{1}{2}} \\
 &\leq (C(1 + \text{MSE}_p(u, v)))^{\frac{1}{2}} \text{MSE}_p(u, v)^{\frac{1}{2}}.
 \end{aligned} \tag{62}$$

703

□

704 C SBI Benchmark

705 In this section, we collect missing details and additional results for the analysis of the SBI benchmark
 706 in Section 4.

707 C.1 Network architecture and hyperparameters

708 For each task and simulation budget in the benchmark, we perform a mild hyperparameter optimization.
 709 We sweep over the batch size and learning rate (which is particularly important as the simulation
 710 budgets differ by orders of magnitudes), the network size and the α parameter for the time prior
 711 defined in Section 3.3 (see Tab. 2 for the specific values). We reserve 5% of the simulation budget for
 712 validation and choose the model with the best validation loss across all configurations.

713 C.2 Additional results

714 We here provide various additional results for the SBI benchmark. First, we compare the performance
 715 of FMPE and NPE when using the Maximum Mean Discrepancy metric (MMD). The results can
 716 be found in Fig. 6. FMPE shows superior performance to NPE for most tasks and simulation
 717 budgets. Compared to the C2ST scores in Fig. 3 the improvement shown by FMPE in MMD is more
 718 substantial.

719 Fig. 7 compares the FMPE results with the optimal transport path from the main text with a comparable
 720 score matching model using the Variance Preserving diffusion path [15]. The score matching results
 721 were obtained using the same batch size, network size and learning rate as the FMPE network, while
 722 optimizing for $\beta_{\min} \in \{0.1, 1, 4\}$ and $\beta_{\max} \in \{4, 7, 10\}$. FMPE with the optimal transport path
 723 clearly outperforms the score-based model on almost all configurations.

724 Finally we compare FMPE using the architecture proposed in Section 3.2 with (t, θ) -conditioning via
 725 gated linear units to FMPE with a naive architecture operating directly on the concatenated (t, θ, x)
 726 vector, see Fig. 8. For the two displayed tasks the context dimension $\dim(x) = 100$ is much larger
 727 than the parameter dimension $\dim(\theta) \in \{5, 10\}$, and there is a clear performance gain in using the
 728 GLU conditioning. Our interpretation is that the low dimensionality of (t, θ) means that it is not
 729 well-learned by the network when simply concatenated with x .

730 Fig. 9 displays the densities of the reference samples under the FMPE model as a histogram for all
 731 tasks (extended version of Fig. 4). The support of the learned model $q(\theta|x)$ covers the reference
 732 samples $\theta \sim p(\theta|x)$, providing additional empirical evidence for the mass-covering behavior theoret-
 733 ically explored in Thm. 1. However, samples from the true posterior distribution may have a small
 734 density under the learned model, especially if the deviation between model and reference is high; see
 735 Lotka-Volterra (bottom right panel).

736 D Gravitational-wave inference

737 We here provide the missing details and additional results for the gravitational wave inference problem
 738 analyzed in Section 5.

hyperparameter	sweep values
hidden dimensions	2^n for $n \in \{4, \dots, 10\}$
number of blocks	$10, \dots, 18$
batch size	2^n for $n \in \{2, \dots, 9\}$
learning rate	$1.e-3, 5.e-4, 2.e-4, 1.e-4$
α (for time prior)	$-0.25, -0.5, 0, 1, 4$

Table 2: Sweep values for the hyperparameters for the SBI benchmark. We split the configurations according to simulation budgets, e.g. for 1000 simulations, we only swept over smaller values for network size and batch size. The network architecture has a diamond shape, with increasing layer width from smallest to largest and then decreasing to the output dimension. Each block consists of two fully-connected residual layers.

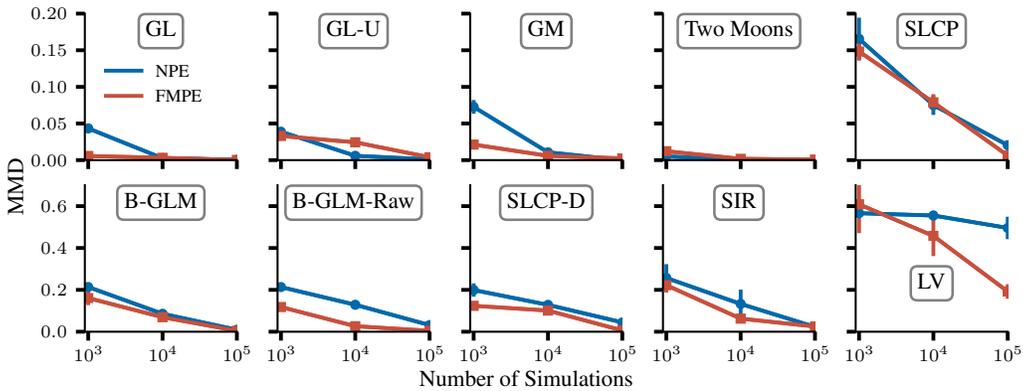


Figure 6: Comparison of FMPE and NPE performance across 10 SBI benchmarking tasks [45]. We here quantify the deviation in terms of the Maximum Mean Discrepancy (MMD) as an alternative metric to the C2ST score used in Fig. 3. MMD can be sensitive to its hyperparameters [45], so we use the C2ST score as a primary performance metric.

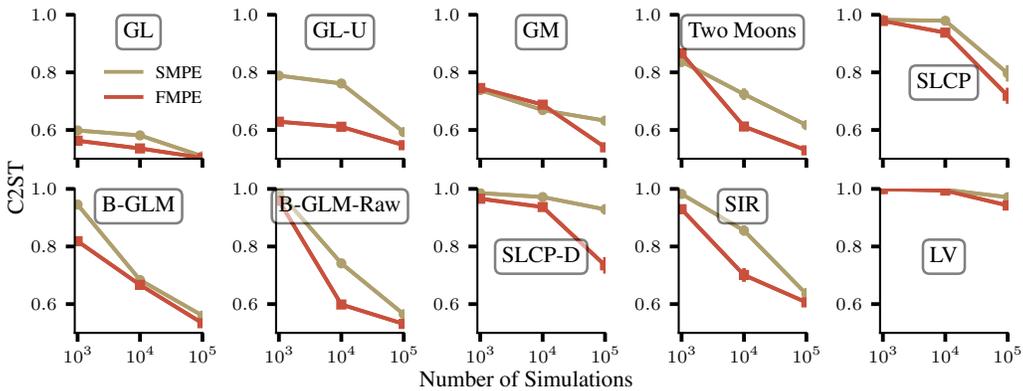


Figure 7: Comparison of FMPE with the optimal transport path (as used throughout the main paper) with comparable models trained with a Variance Preserving diffusion path [15] by regressing on the score (SMPE). Note that the SMPE baseline shown here is not directly comparable to NPSE [8, 9], as this method uses Langevin steps, which reduces the dependence of the results on the vector field for small t (at the cost of a tractable density).

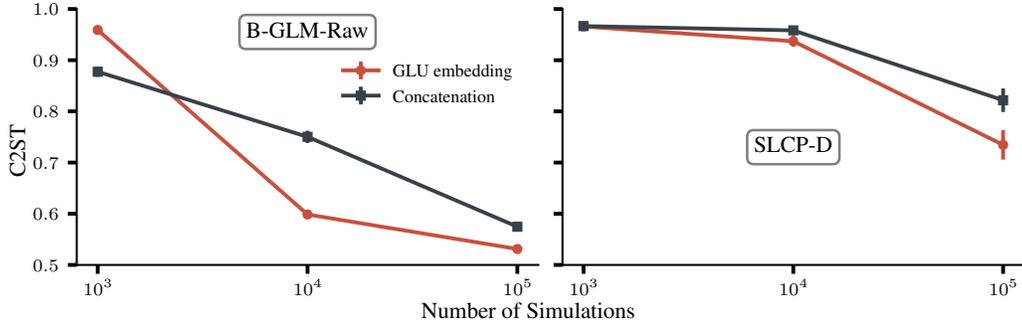


Figure 8: Comparison of the architecture proposed in Section 3.2 with gated linear units for the (t, θ) -conditioning (red) and a naive architecture based on a simple concatenation of (t, θ, x) (black). FMPE with the proposed architecture performs substantially better.

hyperparameter	values
residual blocks	2048, 4096×3 , 2048×3 , 1024×6 , 512×8 , 256×10 , 128×5 , 64×3 , 32×3 , 16×3
residual blocks (t, θ) embedding	16, 32, 64, 128, 256
batch size	4096
learning rate	$5.e-4$
α (for time prior)	1
residual blocks	2048×2 , 1024×4 , 512×4 , 256×4 , 128×4 , 64×3 , 32×3 , 16×3
residual blocks (t, θ) embedding	16, 32, 64, 128, 256
batch size	4096
learning rate	$5.e-4$
α (for time prior)	1

Table 3: Hyperparameters for the FMPE models used in the main text (top) and in the ablation study (bottom, see Fig. 10). The network is composed of a sequence of residual blocks, each consisting of two fully-connected hidden layers, with a linear layer between each pair of blocks. The ablation network is the same as the embedding network that feeds into the NPE normalizing flow.

739 D.1 Network architecture and hyperparameters

740 Compared to NPE with normalizing flows, FMPE allows for generally simpler architectures, since
 741 the output of the network is simply a vector field. This also holds for NPSE (model also defined by a
 742 vector) and NRE (defined by a scalar). Our FMPE architecture builds on the embedding network
 743 developed in [7], however we extend the network capacity by adding more residual blocks (Tab. 3,
 744 top panel). For the (t, θ) -conditioning we use gated linear units applied to each residual block, as
 745 described in Section 3.2. We also use a small residual network to embed (t, θ) before applying the
 746 gated linear units.

747 In this Appendix we also perform an ablation study, using the *same* embedding network as the
 748 NPE network (Tab. 3, bottom panel). For this configuration, we additionally study the effect of
 749 conditioning on (t, θ) starting from different layers of the main residual network.

750 D.2 Data settings

751 We use the data settings described in [7], with a few minor modifications. In particular, we use the
 752 waveform model IMRPhenomPv2 [75–77] and the prior displayed in Tab. 4. Compared to [7], we
 753 reduce the frequency range from $[20, 1024]$ Hz to $[20, 512]$ Hz to reduce the computational load for
 754 data preprocessing. We also omit the conditioning on the detector noise power spectral density (PSD)
 755 introduced in [7] as we evaluate on a single GW event. Preliminary tests show that the performance

Description	Parameter	Prior
component masses	m_1, m_2	$[10, 120] M_\odot, m_1 \geq m_2$
chirp mass	$M_c = (m_1 m_2)^{\frac{3}{5}} / (m_1 + m_2)^{\frac{1}{5}}$	$[20, 120] M_\odot$ (constraint)
mass ratio	$q = m_2 / m_1$	$[0.125, 1.0]$ (constraint)
spin magnitudes	a_1, a_2	$[0, 0.99]$
spin angles	$\theta_1, \theta_2, \phi_{12}, \phi_{JL}$	standard as in [78]
time of coalescence	t_c	$[-0.03, 0.03]$ s
luminosity distance	d_L	$[100, 1000]$ Mpc
reference phase	ϕ_c	$[0, 2\pi]$
inclination	θ_{JN}	$[0, \pi]$ uniform in sine
polarization	ψ	$[0, \pi]$
sky position	α, β	uniform over sky

Table 4: Priors for the astrophysical binary black hole parameters. Priors are uniform over the specified range unless indicated otherwise. Our models infer the mass parameters in the basis (M_c, q) and marginalize over the phase parameter ϕ_c .

756 with PSD conditioning is similar to the results reported in this paper. All changes to the data settings
757 have been applied to FMPE and the NPE baselines alike to enable a fair comparison.

758 D.3 Additional results

759 Tab. 5 displays the inference times for FMPE and NPE. NPE requires only a single network pass
760 to produce samples and (log-)probabilities, whereas many forwards passes are needed for FMPE
761 to solve the ODE with a specific level of accuracy. A significant portion of the additional time
762 required for calculating (log-)probabilities in conjunction with the samples is spent on computing the
763 divergence of the vector field, see Eq. (3).

764 Fig. 10 presents a comparison of the FMPE performance using networks of the same hidden dimen-
765 sions as the NPE embedding network (Tab. 3 bottom panel). This comparison includes an ablation
766 study on the timing of the (t, θ) GLU-conditioning. In the top-row network, the (t, θ) conditioning is
767 applied only after the 256-dimensional blocks. In contrast, the middle-row network receives (t, θ)
768 immediately after the initial residual block. With FMPE we can achieve performance comparable to
769 NPE, while having only $\approx 1/3$ of the network size (most of the NPE network parameters are in the
770 flow). This suggests that parameterizing the target distribution in terms of a vector field requires less
771 learning capacity, compared to directly learning its density. Delaying the (t, θ) conditioning until
772 the final layers impairs performance. However, the number of FLOPs at inference is considerably
773 reduced, as the context embedding can be cached and a network pass only involves the few layers
774 with the (t, θ) conditioning. Consequently, there’s a trade-off between accuracy and inference speed,
775 which we will explore in a greater scope in future work.

	Network Passes	Inference Time (per batch)
FMPE (sample only)	248	26s
FMPE (sample and log probs)	350	352s
NPE (sample and log probs)	1	1.5s

Table 5: Inference times per batch for FMPE and NPE on a single Nvidia A100 GPU, using the training batch size of 4096. We solve the ODE for FMPE using the `dopri5` discretization [79] with absolute and relative tolerances of $1e-7$. For FMPE, generation of the (log-)probabilities additionally requires the computation of the divergence, see equation (3). This needs additional memory and therefore limits the maximum batch size that can be used at inference.

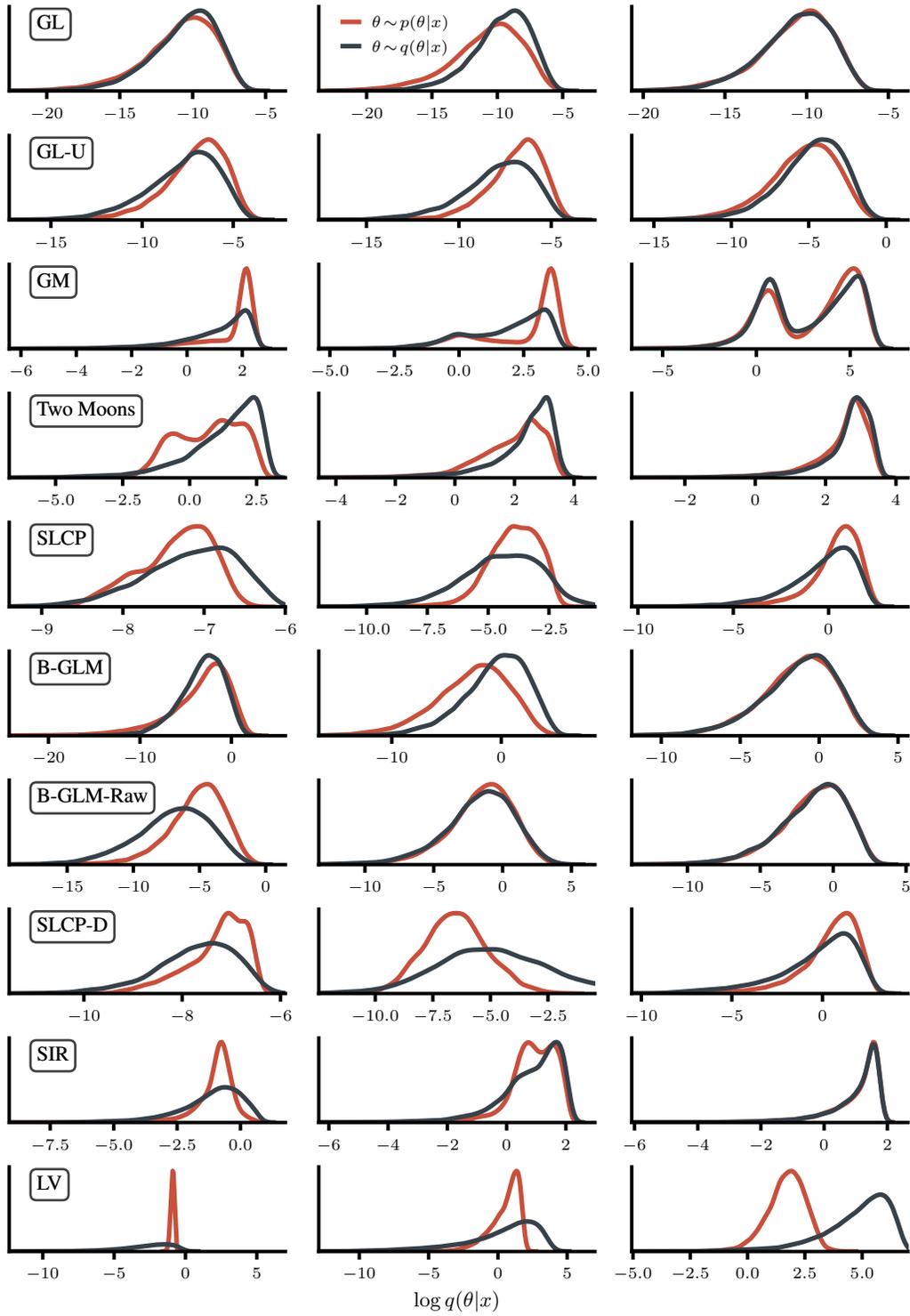


Figure 9: Histogram of FMPE densities $\log q(\theta|x)$ for samples $\theta \sim q(\theta|x)$ and reference samples $\theta \sim p(\theta|x)$ for simulation budgets $N = 10^3$ (left), $N = 10^4$ (center) and $N = 10^5$ (right). The reference samples $\theta \sim p(\theta|x)$ are all within the support of the learned model $q(\theta|x)$, indicating mass covering FMPE results. Nonetheless, reference samples may have a small density under $q(\theta|x)$, if the validation loss is high, see Lotka-Volterra (LV).

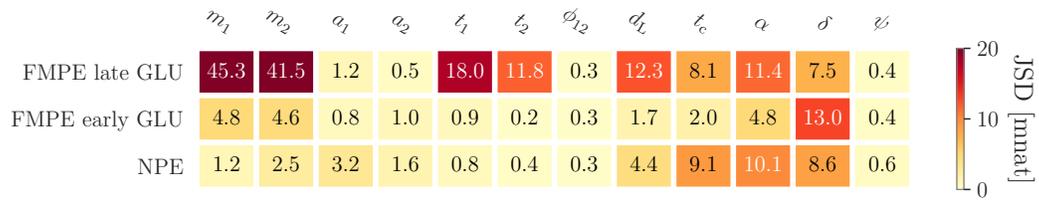


Figure 10: Jensen-Shannon divergence between inferred posteriors and the reference posteriors for GW150914 [72]. We compare two FMPE models with the same architecture as the NPE embedding network, see Tab. 3 bottom panel. For the model in the first row, the GLU conditioning of (θ, t) is only applied before the final 128-dim blocks. The model in the middle row is given the context after the very first 2048 block.