

Automatic Posterior Transformation for Likelihood-free Inference

David S. Greenberg¹ Marcel Nonnenmacher¹ Jakob H. Macke¹

Abstract

How can one perform Bayesian inference on stochastic simulators with intractable likelihoods? A recent approach is to learn the posterior from adaptively proposed simulations using neural network-based conditional density estimators. However, existing methods are limited to a narrow range of proposal distributions or require importance weighting that can limit performance in practice. Here we present automatic posterior transformation (APT), a new sequential neural posterior estimation method for simulation-based inference. APT can modify the posterior estimate using arbitrary, dynamically updated proposals, and is compatible with powerful flow-based density estimators. It is more flexible, scalable and efficient than previous simulation-based inference techniques. APT can operate directly on high-dimensional time series and image data, opening up new applications for likelihood-free inference.

1. Introduction

Many applications in science, engineering and economics make extensive use of complex simulations describing the structure and dynamics of the process being investigated. Such models are derived from knowledge of the mechanisms and principles underlying the data-generating process, and are of critical importance for scientific hypothesis-building and testing. However, linking complex mechanistic models to empirical measurements can be challenging, since many models are defined implicitly through stochastic simulators (e.g. metabolic models, spiking networks in neuroscience, climate and weather simulations, chemical reaction systems, detector models in high energy physics, graphics engines, population models in genetics, dynamical systems in ecology, complex economic models, . . .).

¹Computational Neuroengineering, Department of Electrical and Computer Engineering, Technical University of Munich, Munich, Germany. Correspondence to: <{david.greenberg, marcel.Nonnenmacher, macke}@tum.de>.

Proceedings of the 36th International Conference on Machine Learning, Long Beach, California, PMLR 97, 2019. Copyright 2019 by the author(s).

For complex data-generating processes such as simulation-based models, it is often impossible to compute the likelihood $p(x|\theta)$ of data x given parameters θ , because it involves intractable integrals, because a simulator’s internal states are unavailable or because real-world experiments are involved. Classical approaches to such likelihood-free statistical inference (also known as Approximate Bayesian Computation (ABC), see [Sisson et al., 2018](#)) do not scale to high-dimensional applications, and typically rely on ad-hoc choices to design summary statistics and distance functions.

Recently, several studies have trained conditional density estimators to perform simulation-based inference, while adaptively tuning the simulations to yield informative data. The techniques fall into two main classes, which seek to directly estimate either the likelihood or the posterior:

Synthetic likelihood (SL) methods ([Wood, 2010](#); [Fan et al., 2013](#); [Turner & Sederberg, 2014](#)) aim to estimate the likelihood $p(x|\theta)$, and plug the results into an inference procedure (such as MCMC) to compute the posterior. A powerful recent approach, sequential neural likelihood (SNL), trains a neural conditional density estimator (e.g. masked flow, [Papamakarios et al., 2017](#)) to estimate $p(x|\theta)$ for all x and θ ([Papamakarios et al., 2018](#); [Lueckmann et al., 2018](#)).

Posterior density estimation approaches directly target the posterior $p(\theta|x)$ by training a density-estimation neural network from (simulated) data x to θ . This approach does not require additional inference procedures, and thus naturally *amortizes* inference. In addition, it leverages the ability of neural networks to learn informative features from data.

However, posterior estimation faces a key challenge: Drawing simulation parameters from the prior is wasteful, but other, adaptively-chosen proposals require either numerically unstable post-hoc corrections ([Papamakarios & Murray, 2016](#)) or importance weights ([Lueckmann et al., 2017](#)) that increase variance during learning (details in 2.3).

Here we propose Automatic Posterior Transformation (APT), a new sequential neural posterior estimation approach that combines desirable properties of posterior estimation (directly targeting the posterior, amortization, feature learning) and likelihood estimation (flexible proposals, no importance weights or post-hoc corrections). APT learns a mapping from data to the *true* posterior by maximizing the probability of simulation parameters under the *proposal*

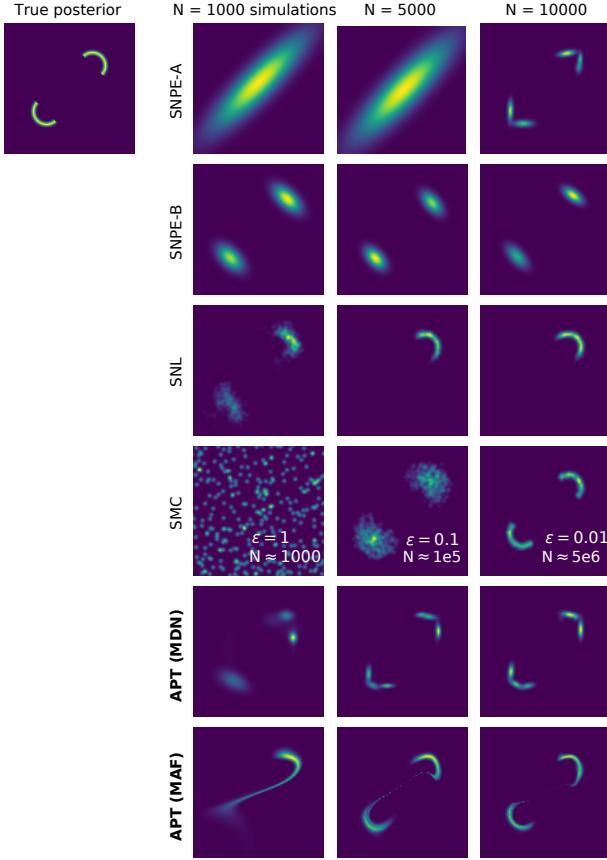


Figure 1. Comparison of inference algorithms on ‘two-moons’ simulator. SNPE-A uses Gaussian posteriors/proposals until the last round. SNPE-B’s importance weights can lead to slow learning. SNL requires MCMC sampling. SMC-ABC requires a distance function, and far more simulations than other methods. APT with mixture- or flow-based density estimators avoids these limitations.

posterior. By recasting inference as a density ratio estimation problem, it can incorporate powerful flow-based density estimators and use arbitrary, dynamically updated proposals.

We demonstrate the effectiveness and flexibility of APT on a variety of problems. It outperforms previous posterior density estimation methods and scales to high dimensional data, with efficient inference on Lokta-Volterra time series and $10k$ -dimensional image data without summary statistics.

2. Simulation-based inference with conditional density estimators

2.1. Problem statement

Suppose a simulator generates synthetic data $x \in \mathbb{R}^d$ for any parameter $\theta \in \mathbb{R}^n$, but the density $p(x|\theta)$ is unknown or intractable. Data x_o are observed from the simulator (or a real-world process it describes) and a prior $p(\theta)$ is known. *Likelihood-free inference* (LFI) seeks to accurately estimate $p(\theta|x_o)$ without access to the likelihood $p(x|\theta)$, by carrying

out a limited number of simulations.

2.2. Conditional density estimation

Likelihood-free inference can be viewed as conditional density estimation, where the task is to map each simulator result x onto an estimate of the posterior density $p(\theta|x)$. Once the density estimator has been trained on simulated data x , it can then be applied to empirical data x_o to compute the posterior. The posterior estimate is selected from a family of densities q_ψ , where ψ are distribution parameters. The mapping from x onto ψ is learned by adjusting the weights ϕ of a neural network F (or another flexible function approximator) so that $q_{F(x,\phi)}(\theta) \approx p(\theta|x)$.

To train the network, we can simulate using prior-drawn parameters to build a dataset $\{(\theta_j, x_j)\}$ and minimize the loss $\mathcal{L}(\phi) = -\sum_{j=1}^N \log q_{F(x_j, \phi)}(\theta_j)$ over network weights ϕ . For sufficiently expressive F and q_ψ , the mapping from x to the posterior $p(\theta|x)$ will be learned as $N \rightarrow \infty$. After training, we estimate the target posterior $p(\theta|x_o)$ by $q_{F(x_o, \phi)}(\theta)$.

2.3. Sequential neural density estimation

Since we are ultimately interested in the posterior at x_o , simulations from parameters with very low posterior density $p(\theta|x_o)$ may not be useful for learning ϕ . Thus, after initially estimating $\theta|x_o$ using simulations from the prior $p(\theta)$, we want future simulations to use a proposal $\tilde{p}(\theta)$ which is more informative about $\theta|x_o$ (Papamakarios & Murray, 2016; Gutmann & Corander, 2016; Lueckmann et al., 2017; Sisson et al., 2007; Blum & François, 2010). This iterative refinement of the posterior estimate and proposal is known as sequential neural posterior estimation (SNPE).

Unfortunately, minimizing \mathcal{L} on samples drawn from a proposal $\tilde{p}(\theta)$ no longer yields the target posterior but rather¹

$$\tilde{p}(\theta|x) = p(\theta|x) \frac{\tilde{p}(\theta) p(x)}{\tilde{p}(\theta) \tilde{p}(x)} \quad (1)$$

where $\tilde{p}(x) = \int_\theta \tilde{p}(\theta)p(x|\theta)$. We call $\tilde{p}(\theta|x)$ the *proposal posterior*. It would be the correct posterior if $\tilde{p}(\theta)$ were the prior. LFI methods employing neural conditional density estimators differ primarily in how they deal with this problem, with three main approaches developed so far:

SNPE-A (Papamakarios & Murray, 2016) trains F to target the proposal posterior $\tilde{p}(\theta|x)$ and then post-hoc solves (1) for $p(\theta|x_o)$. To ensure a closed-form solution, q_ψ is restricted to be a mixture of Gaussians (MoG), $\tilde{p}(\theta)$ to be Gaussian and $p(\theta)$ to be Gaussian or uniform. This approach is simple and can be highly effective, but does not admit multimodal proposals (Fig. 1, first row, details in A.5.1). Furthermore, SNPE-A can return non-positive-

¹We assume $\tilde{p}(\theta) = 0$ where $p(\theta) = 0$.

definite Gaussian covariance matrices when solving (1).

SNPE-B (Lueckmann et al., 2017) minimizes an importance-weighted loss $-\sum_{j=1}^N \frac{p(\theta_j)}{\tilde{p}(\theta_j)} \log q_{F(x_j, \phi)}(\theta_j)$. This allows direct recovery of $p(\theta|x)$ from $q_{F(x, \phi)}$ with no correction and no restrictions on $p(\theta)$, $\tilde{p}(\theta)$ or $q_\psi(\theta)$. However, the importance weights $p(\theta_j)/\tilde{p}(\theta_j)$ greatly increase the variance of parameter updates during learning, which can lead to slow or inaccurate inference (Fig. 1, second row).

SNL instead learns a neural conditional density estimate of the likelihood $p(x|\theta)$, allowing simulation parameters to be drawn from any proposal (e.g. the posterior, Papamakarios et al., 2018) or chosen in an active learning scheme (Lueckmann et al., 2018). However, instead of directly inferring $p(\theta|x_o)$, it uses additional MCMC sampling that can be costly or inefficient for complex posteriors (Fig. 1, third row). Estimating the likelihood (rather than the posterior) can also be more difficult in some cases (see below). Nonetheless, SNL is accurate and efficient on many problems (Durkan et al., 2018), and generally outperforms classical ABC approaches such as SMC-ABC.

3. Automatic Posterior Transformation

APT² is a SNPE technique which combines desirable properties of existing posterior density estimation and likelihood-based approaches. APT learns to infer the *true* posterior by maximizing an estimated *proposal* posterior. It uses (1) to form a parameterization which makes it possible to automatically transform between estimates of $p(\theta|x)$ and $\tilde{p}(\theta|x)$, and thus to easily ‘read-off’ the posterior estimate.³ We will show how this trick avoids the numerical challenges of previous SNPE techniques, and makes it possible to use a wide range of proposals and density estimators.

APT uses $q_{F(x, \phi)}(\theta)$ to represent an estimate of $p(\theta|x)$. To transform this into an estimate of $\tilde{p}(\theta|x)$, we first observe that by (1), $\tilde{p}(\theta|x) \propto p(\theta|x)\tilde{p}(\theta)/p(\theta)$. We therefore define

$$\tilde{q}_{x, \phi}(\theta) = q_{F(x, \phi)}(\theta) \frac{\tilde{p}(\theta)}{p(\theta)} \frac{1}{Z(x, \phi)}, \quad (2)$$

$Z(x, \phi) = \int_\theta q_{F(x, \phi)}(\theta) \frac{\tilde{p}(\theta)}{p(\theta)}$ is a normalization constant.

We train the network to carry out posterior inference by minimizing $\tilde{\mathcal{L}}(\phi) = -\sum_{j=1}^N \log \tilde{q}_{x, \phi}(\theta_j)$. This procedure recovers both the true and proposal posteriors, as follows:

If the conditional density estimator $q_{F(x, \phi)}$ is expressive enough that some ϕ^* exists for which $q_{F(x, \phi^*)}(\theta) = p(\theta|x)$, then by (1) also $\tilde{q}_{x, \phi^*}(\theta) = \tilde{p}(\theta|x)$. Therefore by Prop. 1 of (Papamakarios & Murray, 2016), minimizing $\tilde{\mathcal{L}}(\phi)$ yields $q_{F(x, \phi)}(\theta) \rightarrow p(\theta|x)$ and $\tilde{q}_{x, \phi}(\theta) \rightarrow \tilde{p}(\theta|x)$ as $N \rightarrow \infty$.

²Or SNPE-C in the taxonomy of (Papamakarios et al., 2018).

³Code available at github.com/mackelab/delfi.

Similar to SNPE-A/B and SNL, APT iteratively refines the network weights ϕ and proposal $\tilde{p}(\theta)$ over multiple simulation rounds. Each round uses a different proposal $\tilde{p}_r(\theta)$, leading to different transformations in (2) when calculating the loss $\tilde{\mathcal{L}}$. However, since $q_{F(x, \phi)}(\theta) = p(\theta|x)$ minimizes the expectation $\mathbb{E}[\tilde{\mathcal{L}}]$ for *any* proposal, APT can train on data from multiple rounds simply by adding their loss terms together. In contrast, SNPE-A cannot re-use data across rounds and SNPE-B must apply different importance weights. Algorithm 1 describes APT in the case that each round’s final posterior estimate is the next round’s proposal.

APT supports a wide range of proposals and density estimators, including mixture-density networks and flows (Fig. 1, last two rows). Its only requirement is a closed-form solution of (2) for $\tilde{q}_{F(x, \phi)}(\theta)$, to be optimized during learning and sampled from afterwards. We next describe several combinations of $p(\theta)$, $\tilde{p}(\theta)$ and $q_\psi(\theta)$ for which this is possible.

Algorithm 1 APT with per-round proposal updates

Input: simulator with (implicit) density $p(x|\theta)$, data x_o , prior $p(\theta)$, density family q_ψ , neural network $F(x, \phi)$, simulations per round N , number of rounds R .

```

 $\tilde{p}_1(\theta) := p(\theta)$ 
for  $r = 1$  to  $R$  do
    for  $j = 1$  to  $N$  do
        Sample  $\theta_{r,j} \sim \tilde{p}_r(\theta)$ 
        Simulate  $x_{r,j} \sim p(x|\theta_{r,j})$ 
    end for
     $\phi \leftarrow \operatorname{argmin}_\phi \sum_{i=1}^r \sum_{j=1}^N -\log \tilde{q}_{x_{i,j}, \phi}(\theta_{i,j})$  using (2)
     $\tilde{p}_{r+1}(\theta) := q_{F(x_o, \phi)}(\theta)$ 
end for
return  $q_{F(x_o, \phi)}(\theta)$ 

```

3.1. Gaussian and Mixture-of-Gaussians proposals

In the simplest case all distributions are Gaussian, with prior $\mathcal{N}_\theta(\mu_0, S_0)$, proposal $\mathcal{N}_\theta(\tilde{\mu}_0, \tilde{S}_0)$, posterior $\mathcal{N}_\theta(\mu, S)$ and proposal posterior $\mathcal{N}_\theta(\tilde{\mu}, \tilde{S})$. Then (1-2) reduce to

$$\tilde{S}^{-1} = S^{-1} + \tilde{S}_0^{-1} - S_0^{-1} \quad (3)$$

$$\tilde{S}^{-1}\tilde{\mu} = S^{-1}\mu + \tilde{S}_0^{-1}\tilde{\mu}_0 - S_0^{-1}\mu_0 \quad (4)$$

While SNPE-A uses these relations to calculate the true posterior from the proposal posterior *after* learning, APT does the opposite *during* learning. Since the proposal is narrower than the prior in all sensible scenarios, $\tilde{S}_0^{-1} - S_0^{-1}$ is positive definite and APT cannot produce invalid \tilde{S} .

The updates (3-4) can readily be extended to uniform priors by setting $S_0^{-1} = 0$, and to MoG posteriors by sum-

ming over components (Papamakarios & Murray, 2016). When $\tilde{p}(\theta)$ is a MoG with L components and $q_{F(x,\phi)}(\theta)$ a MoG with K components, solving (2) for $\tilde{q}_{x,\phi}(\theta)$ yields an LK -component MoG (A.1). This allows APT to propose simulation parameters from multimodal distributions.

Table 1. Properties of posterior inference techniques

ALGORITHM	$\tilde{p}(\theta)$	$p(\theta)$	q_ψ
SMC-ABC	ANY	ANY	DISCRETE
SNPE-A	GAUSS	GAUSS/UNI	MDN
SNPE-B	ANY	ANY	ANY
SNL	ANY	ANY	ANY (MCMC)
APT	ANY	ANY	ANY

3.2. Atomic proposals

We would like to extend APT to arbitrary choices of the density estimator, proposal and prior, and especially to powerful flow-based density estimators (Rezende & Mohamed, 2015; Papamakarios et al., 2017; Kingma & Dhariwal, 2018). However, in many cases we cannot solve the integral defining Z in (2). In this section, we show how APT can be trained using ‘atomic’ proposals that only consider a finite set of parameter vectors (‘atoms’) for each simulation, replacing integrals by sums. Provided that each simulation’s parameters are drawn from a different atomic proposal, and that the overall range of possible atoms covers the posterior support, we can infer the full, continuous posterior.

We set $\tilde{p}(\theta) = U_\Theta$, where U_Θ is uniform on $\Theta = \{\theta_1, \dots, \theta_M\}$. Then $\tilde{p}(\theta|x)$ and $\tilde{q}_{x,\phi}(\theta)$ are categorical distributions, $\tilde{\mathcal{L}}$ is a cross-entropy loss and (1-2) reduce to

$$\tilde{p}(\theta|x) = \frac{p(\theta|x)/p(\theta)}{\sum_{\theta' \in \Theta} p(\theta'|x)/p(\theta')} \quad (5)$$

$$\tilde{q}_{x,\phi}(\theta) = \frac{q_{F(x,\phi)}(\theta)/p(\theta)}{\sum_{\theta' \in \Theta} q_{F(x,\phi)}(\theta')/p(\theta')} \quad (6)$$

For fixed Θ , $\mathbb{E}_{\theta \sim U_\Theta, x \sim p(x|\theta)}[\tilde{\mathcal{L}}]$ is minimized precisely when $\tilde{q}_{x,\phi}(\theta) = \tilde{p}(\theta|x)$. In that case $\forall \theta_1, \theta_2 \in \Theta$, $\frac{q_{F(x,\phi)}(\theta_1)}{q_{F(x,\phi)}(\theta_2)} = \frac{p(\theta_1|x)}{p(\theta_2|x)}$ by (5-6)—that is, posterior density ratios are correct over all atoms and for every x the atoms can generate.

Now suppose that Θ is itself sampled from a ‘hyperproposal’ $V(\Theta)$ for each simulation. To minimize $\mathbb{E}[\tilde{\mathcal{L}}]$, the network must then infer the correct posterior density ratios for every pair of samples from each Θ . We formalize this argument as follows (proof in A.4):

Proposition 1. Let $\Theta \sim V$, and let

$$\rho(x, \Theta, \phi) = \begin{cases} 1, & \text{if } \frac{p(\theta_1|x)}{p(\theta_2|x)} = \frac{q_{F(x,\phi)}(\theta_1)}{q_{F(x,\phi)}(\theta_2)} \quad \forall \theta_1, \theta_2 \in \Theta \\ 0, & \text{otherwise} \end{cases}$$

and suppose that for some ϕ^* , $q_{F(x,\phi^*)} = p(\theta|x)$. Then for any ϕ with $\mathbb{E}_{\Theta \sim V, \theta \sim U_\Theta, x \sim p(x|\theta)}[\tilde{\mathcal{L}}(\phi) - \tilde{\mathcal{L}}(\phi^*)] = 0$, $\mathbb{E}_{\Theta \sim V, \theta \sim U_\Theta, x \sim p(x|\theta)}[\rho(x, \Theta, \phi)] = 1$.

By Prop. 1, APT recovers the full posterior shape as $N \rightarrow \infty$, provided that each Θ is constructed by sampling θ ’s independently from a distribution covering the support of $p(\theta|x_o)$. Furthermore, we can get this consistency guarantee using a fixed, nonzero fraction of the atomic proposals to cover the posterior support, while arbitrarily assigning the rest (e.g. for active learning).

Atomic APT’s training data consists of triples (θ, x, Θ) . $q_{F(x,\phi)}$ is plugged into (6) to yield $\tilde{q}_{x,\phi}$, which is evaluated on θ to calculate $\tilde{\mathcal{L}}$. In practice, instead of sampling Θ , θ and x during training we simulate once per round as in non-atomic APT, then resample the triples (θ, x, Θ) from all rounds’ simulations during training. We describe the full algorithm and its computational complexity in A.2.

Learning with such ‘atomic’ proposals has an intuitive interpretation: we are training the network to solve multiple choice test problems, of the format “which of these θ ’s generated this x ?” Given a list of M possible alternatives, the network gives a Bayesian answer by assigning mass to each $\theta \in \Theta$. Prop. 1 shows we can learn to infer continuous posteriors by training on multiple choice questions.

By replacing integrals with sums, we can use flows or any other distributions for the prior, proposal and posterior estimate, so long as the densities can be easily evaluated, we can sample from the proposal and the posterior is differentiable in ψ . Atomic proposals also greatly enhance flexibility when generating new simulations: they can emphasize the posterior’s peaks or tails, or use active learning schemes (Lueckmann et al., 2018; Järvenpää et al., 2018) incompatible with previous SNPE approaches.

3.3. Truncated priors and posteriors

With atomic proposals and a prior with limited support, Prop. 1 only ensures we can retrieve the posterior $p(\theta|x)$ up to an unknown scale factor: $q_{F(x,\phi)}$ can be trained to match its shape, but not its amplitude. In this case we can easily obtain posterior samples by sampling from $q_{F(x,\phi)}$ and rejecting when $p(\theta) = 0$. This is less convenient than directly sampling from the posterior, but still simpler and more efficient than a full MCMC scheme. We discuss the effects of truncated priors on posterior estimation in A.3.

4. Experiments

We compare APT to SNPE-A, SNPE-B and SNL on several problems (implementation details in A.5). (Papamakarios et al., 2018) quantitatively compared these algorithms with classical ABC methods (SMC-ABC and SL).

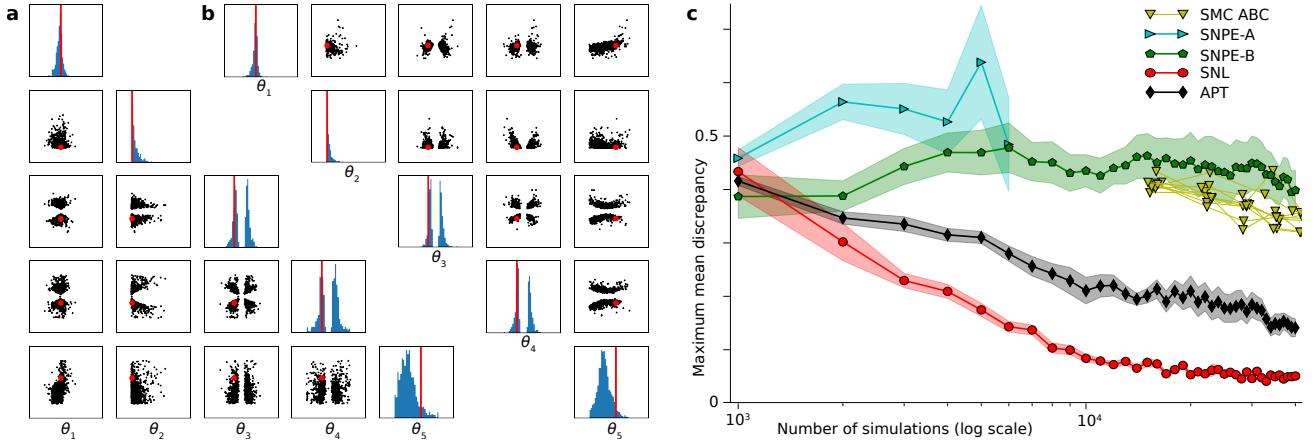


Figure 2. Performance on a model with a simple likelihood and complex posterior (SLCP). **a** Posterior samples from flow-based APT, and **b** SNL (ground-truth parameters in red). **c** Maximum mean discrepancies between estimated and ground-truth posteriors. Mean \pm SEM over 10 random initializations of the simulator and inference methods with identical x_o . SNPE-A terminated after round 6. APT outperforms previous SNPE methods, closing over half the performance gap to SNL on a problem designed to favor likelihood estimation.

4.1. Illustrative toy example with multiple modes

We first illustrate properties and common failure cases of inference methods on a simple ‘two moons’ simulator with 2D parameters θ and observations x (Fig. 1, details in A.5.1, quantitative evaluation in Fig. 7). Its conditionals $p(\theta|x)$ have both global (bimodality) and local (crescent shape) structure. The posterior was learned over 10 rounds of 1000 simulations. We can see clear differences between methods:

SNPE-A is limited to Gaussian proposals and hence cannot learn fine structure until the final round. SNPE-B has MoG proposals, but its importance weights can lead to slow learning: here it fails to capture the crescent shapes. SNL can flexibly approximate the crescent-shaped likelihood $p(x|\theta)$, but then requires an additional inference step to calculate the posterior from the synthetic likelihood—the coordinate-alternating slice sampling method previously used by (Papamakarios et al., 2018) here fails to mix well between the two non-axis-aligned modes (we emphasize that other MCMC approaches or parameter settings might work better here—however, multimodal and high-dimensional posteriors are generally challenging for MCMC). Classical SMC-ABC (Sisson et al., 2007; Beaumont et al., 2009) requires a large number of simulations (details in A.5.1; note higher number of simulations used for SMC).

APT with a Mixture-Density Network (MDN) for density estimation identifies the two-moons structure, and uses it to efficiently guide proposals for subsequent rounds. APT can also be flexibly applied to other density estimators, as illustrated here using a masked autoregressive flow (Papamakarios et al., 2017) to represent the posterior.

4.2. Toy example with simple likelihood but complex posterior

(Papamakarios et al., 2018) introduce a toy example designed to have a Simple Likelihood and Complex Posterior (SLCP model, details in A.5.2). Its likelihood $p(x|\theta)$ is a Gaussian in x whose mean and covariance depend nonlinearly on θ , while its posteriors $p(\theta|x)$ are non-Gaussian and multi-modal in θ . Thus, estimating the posterior $p(\theta|x)$ is more challenging than estimating the likelihood $p(x|\theta)$. Nevertheless, we demonstrate that with flexible flow-based conditional density estimators, APT’s posterior estimates are similar to SNL’s (Fig. 2a vs. b, ground-truth posterior in Papamakarios et al., 2018). A quantitative evaluation (using Maximum Mean Discrepancy, Gretton et al., 2012) shows that APT outperforms SNPE-A/B, and is closer in performance to SNL, on a problem which is designed to favor SNL.

4.3. Effect of non-informative observations

Posterior inference for high-dimensional data is a challenging problem (Ong et al., 2018), especially when multiple data dimensions are uninformative about θ . This scenario is common in neuroscience (Paninski & Cunningham, 2017), physics (Brehmer et al., 2018), systems biology (Clarke et al., 2008) and inverse graphics (Romaszko et al., 2017).

While classical ABC methods require model-specific low-dimensional summary statistics, APT maps from x to a posterior estimate using deep neural networks, which excel at learning relevant features from data. To test whether this would allow efficient inference in models with many uninformative data dimensions, we appended uninformative outputs drawn from a mixture of Student-t distributions to the 8D SLCP outputs (details in A.5.2). These noise outputs

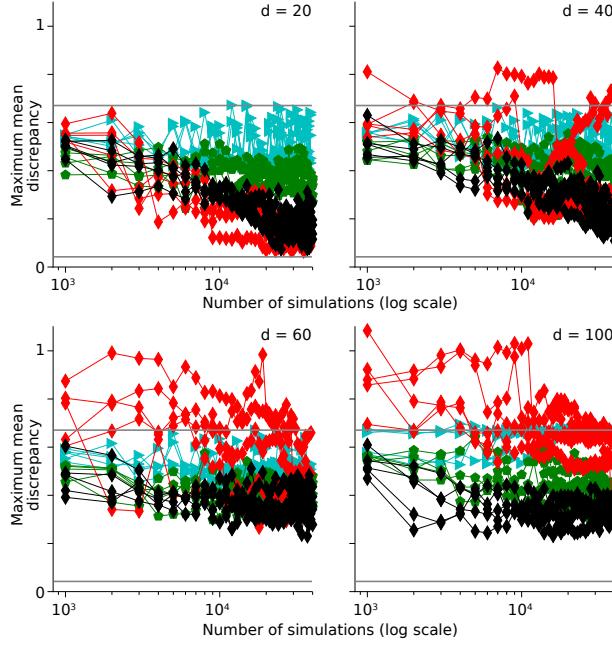


Figure 3. Inference with uninformative data dimensions. Appending uninformative dimensions to the SLCP model demonstrates APT can learn relevant features. Maximum mean discrepancies (MMDs) between estimated and ground-truth posteriors (5 random initializations) for each value of $d = \dim(x)$. Gray lines show MMDs for prior (upper) and ground-truth posterior samples.

are generated independently from informative outputs, so the true posterior $p(\theta|x_o)$ is unaffected but inference becomes more difficult as the data dimensionality d increases.

At moderate d both SNL and APT can recover $p(\theta|x_o)$, but at large d SNL’s posterior estimates are no better than the prior, while APT degrades only slightly (Fig. 3). Presumably, SNL ‘invests’ too many resources into estimating densities on irrelevant dimensions. For fixed $\dim(\theta)$, the number of network weights grows linearly with d for posterior density estimators, but quadratically for MDN- or flow-based likelihood estimators.

4.4. Population ecology model

The Lotka-Volterra model of coupled predator and prey populations (Lotka, 1920) is a classical likelihood-free inference benchmark. $\theta \in \mathbb{R}^4$ governs growth rates and predator-prey interactions, and x consists of population counts (here at 150 fixed intervals). Simulation results are typically compressed into a set of summary statistics \bar{x} with $d = 9$.

This simulator, with complex \bar{x} distributions (Fig. 4b) and simple posteriors (Fig. 4a), is well-suited to posterior density estimation. APT produced tight posteriors around the true parameters with fewer simulations than other methods (Fig. 4c). While the ground truth posterior is unavailable for

this model, for large N APT and SNL compute similar posterior estimates in very different ways, suggesting a close approximation has been obtained (see A.5.3).

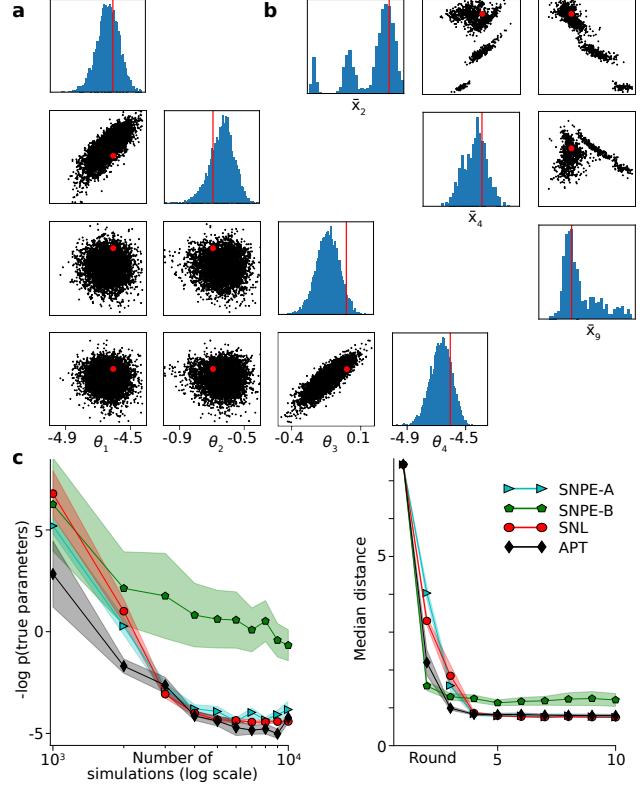


Figure 4. Lotka-Volterra with summary statistics. The Lotka-Volterra model has simple posteriors and complex conditionals $p(\bar{x}|\theta)$. **a** Posterior estimated by APT (ground-truth parameters in red). **b** Samples from $p(\bar{x}_2, \bar{x}_4, \bar{x}_9|\theta)$ for ground-truth parameters (\bar{x}_o in red). **c** Negative log-probabilities of ground-truth parameters and median distances $|\bar{x} - \bar{x}_o|$ under the posterior estimates across rounds (means \pm SEM across 10 different random initializations).

4.5. Inference on raw time-series

A central problem in ABC is the fact that virtually all algorithms rely on summary statistics in order to reduce the dimensionality of the data (Fearnhead & Prangle, 2012; Blum et al., 2013; Jiang et al., 2017)—designing summary statistics can require domain-specific knowledge and/or separate optimization procedures, and might bias algorithms if the statistics do not adequately summarize the informative features of the data. To investigate whether sequential neural posterior estimation with APT can be applied to raw simulator outputs, we also applied APT using a recurrent neural network (RNN) to directly map x onto a posterior density estimate (details in A.5.3). RNN-APT’s estimate of $p(\theta|x)$ (Fig. 5a) was nearly identical to the previous $p(\theta|\bar{x})$. Over repeated runs of the algorithm the posterior marginals remained tightly clustered around the true parameters in both cases (Fig. 5c). This shows that APT can operate

without summary statistics on high-dimensional data, and suggests that $p(\theta|\bar{x})$ may be close to $p(\theta|x)$.

We also examined a more difficult version of the same problem that has been used as an LFI benchmark (Owen et al., 2015), where the observed counts are corrupted by observation noise (Fig. 5d). In this case posterior distributions inferred by the RNN were broader, but still clustered around the true parameters. Posteriors inferred from summary statistics were significantly broader and farther from the true parameters than those of the RNN, suggesting that the summary statistics may discard information about the parameters in this case. By training APT end-to-end we were able to automatically extract relevant information from the data, without using bespoke summary statistics.

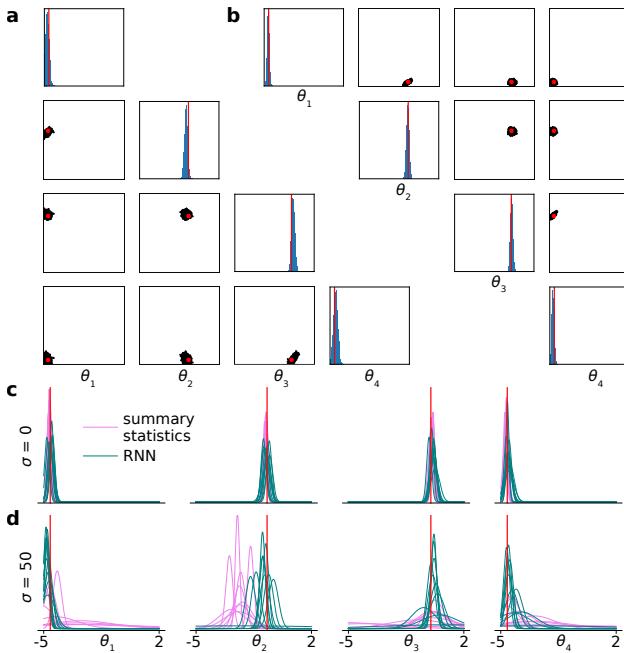


Figure 5. Lotka-Volterra with RNN. **a** Posterior estimate from RNN-based APT (ground-truth parameters in red). **b** Posterior estimate from APT with summary statistics (same data). **c** Estimated posterior marginals from 10 random initializations, using APT with the RNN (teal) or summary statistics (magenta). **d** As in ‘c’ but with i.i.d. Gaussian observation noise, $\sigma = 50$.

4.6. Rock-paper-scissors reaction-diffusion model

Simulators that produce images pose a particular challenge for inference, as they produce high-dimensional data without well-tested summary statistics. In evolutionary game theory, rock-paper-scissors (RPS) models (May & Leonard, 1975) describe nontransitive predator-prey interactions in which species A preys on B, B preys on C and C preys on A, often exhibiting stable biodiversity over long timescales. When each species’ population density varies over a 2-dimensional space, the resulting stochastic partial differen-

tial equation (SPDE) exhibits complex, dynamically shifting spatial structure (Reichenbach et al., 2007). The system’s stability and structure depend on its growth, predation and diffusion rates.

We simulated from the SPDE-RPS model, spatially discretized on a 100×100 lattice and initialized at the unstable uniform steady state. The observation $x_o \in \mathbb{R}^{10000}$ was the final system state, which can be displayed as a 3-channel image (Fig. 6a). To test whether APT could identify relevant features and infer parameter posteriors from high-dimensional image data, we used a convolutional neural network (CNN) as $F(x, \phi)$. APT inferred posteriors that closely encompassed the ground truth parameters used to generate each x_o (Fig. 6c), and posterior-drawn parameters produced simulations that visually resembled x_o (Fig. 6b). Running SNPE-A/B with the same CNN yielded lower posterior probability for the ground truth parameters (details in A.5.4).

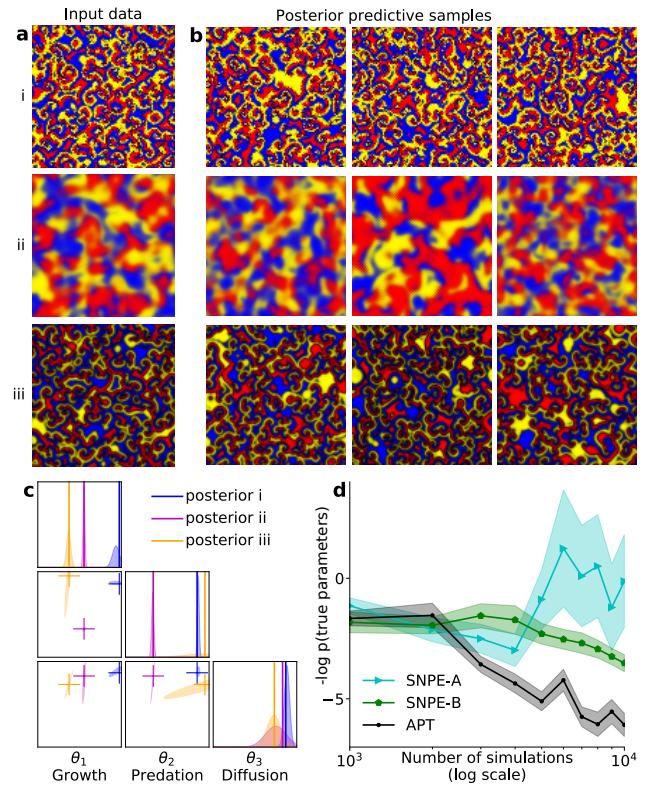


Figure 6. Rock-paper-scissors model with CNN. **a** Observations x_o for 3 RPS parameter sets. **b** 3 APT posterior predictive samples x for each x_o in ‘a’: $\theta \sim q_{F(x_o, \phi)}(\theta)$, $x \sim p(x|\theta)$. **c** APT posteriors for each x_o in ‘a.’ On-diagonal graphs show marginals; off-diagonals show 2 s.d.s around the mean. Solid lines and crosses show ground truth parameters. **d** Negative log-probabilities of ground-truth parameters; mean \pm SEM over 30 inference problems with ground truth parameters drawn from the prior.

4.7. Additional experiments

Finally, we also evaluated two additional benchmark models. On the M/G/1 model (Papamakarios & Murray, 2016), APT outperformed SNPE-A/B and approached SNL in performance (Fig. 9). On the GLM model (Lueckmann et al., 2017) APT was slightly more accurate than SNPE-A, while SNL and SNPE-B were considerably less accurate (Fig. 10).

5. Related work

Classical ABC approaches (Sisson et al., 2018), are based on simulating data from proposal distributions (Pritchard et al., 1999; Beaumont et al., 2002; Marjoram et al., 2003; Beaumont et al., 2009; Gutmann & Corander, 2016) and accepting samples that are sufficiently similar to the observed data. However, this approach generally requires choosing low-dimensional summary statistics, distance functions and rejection thresholds, and is exact only in the limit of high rejection rates. Several papers have investigated procedures for (semi) automatically designing summary statistics (Fearnhead & Prangle, 2012; Blum et al., 2013; Jiang et al., 2017; Dinev & Gutmann, 2018).

Approaches based on estimating the likelihood (Wood, 2010; Fan et al., 2013; Turner & Sederberg, 2014; Papamakarios et al., 2018; Lueckmann et al., 2018) or posterior (Le et al., 2017; Papamakarios & Murray, 2016; Lueckmann et al., 2017; Chan et al., 2018) (which can be traced back to regression-adjustment, Beaumont et al., 2002; Blum & François, 2010) are reviewed above (2.3). Posterior inference with end-to-end feature learning has previously been demonstrated with RNNs (Lueckmann et al., 2017) and exchangeable networks (Chan et al., 2018). Some LFI methods instead aim to infer a point or local estimate of the parameters (Pesah et al., 2018; Louppe & Cranmer, 2017; McCarthy et al., 2017).

Atomic APT trains the posterior model through a series of multiple-choice problems, by minimizing the cross-entropy loss from supervised classification. Several recent studies investigated the use of (neural) classifiers for hypothesis testing or estimation of likelihood ratios: (Dutta et al., 2016; Gutmann et al., 2018) train binary classifiers to discriminate between the conditional and the marginal distribution, (Cranmer et al., 2015) trains classifiers to approximate likelihood ratios for frequentist parameter estimation and hypothesis testing and (Brehmer et al., 2018) presents approaches for learning likelihoods (and scores) from simulator outputs. (Tran et al., 2017a) trains a log ratio estimator to minimize a variational objective, extending neural posterior estimation to implicit density families with intractable $q_\psi(\theta)$, but uses prior-sampled parameters and requires training two networks simultaneously. Atomic APT differs from these in using a single parameterized posterior estimate to derive a

different classifier for each multiple choice problem.

If the simulator also provides access to additional information (e.g. internal variables and their likelihoods or an unbiased estimate of $p(x|\theta)$), this information can be exploited to improve efficiency (Tran et al., 2017a; Brehmer et al., 2018; Andrieu et al., 2010; Tran et al., 2017b). Approaches for so-called *implicit* generative models generally require the simulator to be differentiable with respect to model parameters (Mohamed & Lakshminarayanan, 2016; Huszár, 2017, and references therein). While generative adversarial networks (GANs) also train classifiers on simulated data, our learning objective is different (Huszár, 2017), as in particular as we have no ‘adversary.’ Nevertheless, adversarial training can also be extended to inference problems, and its application to likelihood-free inference has been proposed (Louppe & Cranmer, 2017). It might be possible to combine adversarial learning with APT, especially when using atomic proposals.

6. Discussion

Density estimation is an attractive approach to simulation-based Bayesian inference. We presented APT, which can be applied to arbitrary proposals and a wide range of conditional density estimators. Flexibility in choosing proposals opens up multiple new applications and opportunities for extensions, e.g. combining posterior estimation with active-learning rules, or loss-calibrated likelihood-free inference. APT is more simulation-efficient than previous posterior density estimation approaches on a variety of problems, and can scale to higher dimensional observations.

Depending on the model and analysis problem at hand, the likelihood or the posterior might be easier to approximate. We do not claim that directly targeting the posterior will always outperform synthetic-likelihood approaches. However, it has the advantage of directly yielding a mapping from data to posterior without an additional inference step, and therefore amortizes inference. An advantage of targeting the posterior directly is that specialized neural network architectures can be used to exploit known structure in the data, as we have shown using RNNs and CNNs for time series and image data. Exchangeable neural networks (Zaheer et al., 2017; Chan et al., 2018; Bloem-Reddy & Teh, 2019) could also be applied when x_o consists of multiple i.i.d. observations from the same parameters.

We hope that the combination of flexible conditional density estimators, effective schemes for adaptively choosing simulations, and stable learning frameworks will eventually make statistical inference efficient, easy and (ideally) automated, for a wide range of simulation-based models.

Acknowledgements

We thank Jan-Matthis Lueckmann, Álvaro Tejero-Cantero, Poornima Ramesh, Pedro J. Gonçalves, Jan Böltz, Michael Deistler and Artur Speiser for comments and discussions, Poornima Ramesh for pointing us to the lattice RPS model, Tobias Reichenbach for RPS code and the anonymous reviewers for detailed and constructive feedback. Funded by the German Research Foundation (DFG) through SFB 1233 (276693517), SFB 1089 and SPP 2041 and the German Federal Ministry of Education and Research (BMBF, project ‘ADMIMEM’, FKZ 01IS18052 A-D).

References

- Andrieu, C., Doucet, A., and Holenstein, R. Particle markov chain monte carlo methods. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 72(3): 269–342, 2010.
- Beaumont, M. A., Zhang, W., and Balding, D. J. Approximate bayesian computation in population genetics. *Genetics*, 162(4):2025–2035, 2002.
- Beaumont, M. A., Cornuet, J., Marin, J., and Robert, C. P. Adaptive approximate bayesian computation. *Biometrika*, 2009.
- Bloem-Reddy, B. and Teh, Y. W. Probabilistic symmetry and invariant neural networks. *arXiv preprint arXiv:1901.06082*, 2019.
- Blum, M. G., Nunes, M. A., Prangle, D., Sisson, S. A., et al. A comparative review of dimension reduction methods in approximate bayesian computation. *Statistical Science*, 28(2):189–208, 2013.
- Blum, M. G. B. and François, O. Non-linear regression models for approximate bayesian computation. *Statistics and Computing*, 20(1), 2010.
- Brehmer, J., Louppe, G., Pavéz, J., and Cranmer, K. Mining gold from implicit models to improve likelihood-free inference. *arXiv preprint arXiv:1805.12244*, 2018.
- Chan, J., Perrone, V., Spence, J. P., Jenkins, P. A., Mathieson, S., and Song, Y. S. A likelihood-free inference framework for population genetic data using exchangeable neural networks. *arXiv preprint arXiv:1802.06153*, 2018.
- Cho, K., Van Merriënboer, B., Gulcehre, C., Bahdanau, D., Bougares, F., Schwenk, H., and Bengio, Y. Learning phrase representations using rnn encoder-decoder for statistical machine translation. *arXiv preprint arXiv:1406.1078*, 2014.
- Clarke, R., Ressom, H. W., Wang, A., Xuan, J., Liu, M. C., Gehan, E. A., and Wang, Y. The properties of high-dimensional data spaces: implications for exploring gene and protein expression data. *Nature reviews cancer*, 8(1): 37, 2008.
- Cranmer, K., Pavéz, J., and Louppe, G. Approximating likelihood ratios with calibrated discriminative classifiers. *arXiv preprint arXiv:1506.02169*, 2015.
- Dinev, T. and Gutmann, M. U. Dynamic likelihood-free inference via ratio estimation (dire). *arXiv preprint arXiv:1810.09899*, 2018.
- Durkan, C., Papamakarios, G., and Murray, I. Sequential neural methods for likelihood-free inference. *arXiv preprint arXiv:1811.08723*, 2018.
- Dutta, R., Corander, J., Kaski, S., and Gutmann, M. U. Likelihood-free inference by ratio estimation. *arXiv preprint arXiv:1611.10242*, 2016.
- Fan, Y., Nott, D. J., and Sisson, S. A. Approximate bayesian computation via regression density estimation. *Stat*, 2(1), 2013.
- Fearnhead, P. and Prangle, D. Constructing summary statistics for approximate bayesian computation: semi-automatic approximate bayesian computation. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 74(3):419–474, 2012.
- Gretton, A., Borgwardt, K. M., Rasch, M. J., Schölkopf, B., and Smola, A. A kernel two-sample test. *Journal of Machine Learning Research*, 13(Mar):723–773, 2012.
- Gutmann, M. U. and Corander, J. Bayesian optimization for likelihood-free inference of simulator-based statistical models. *The Journal of Machine Learning Research*, 17 (1):4256–4302, 2016.
- Gutmann, M. U., Dutta, R., Kaski, S., and Corander, J. Likelihood-free inference via classification. *Statistics and Computing*, 28(2):411–425, 2018.
- Huszár, F. Variational inference using implicit distributions. *ArXiv e-prints*, 2017.
- Järvenpää, M., Gutmann, M. U., Pleska, A., Vehtari, A., Marttinen, P., et al. Efficient acquisition rules for model-based approximate bayesian computation. *Bayesian Analysis*, 2018.
- Jiang, B., Wu, T.-y., Zheng, C., and Wong, W. H. Learning summary statistic for approximate bayesian computation via deep neural network. *Statistica Sinica*, pp. 1595–1618, 2017.

- Kingma, D. P. and Dhariwal, P. Glow: Generative flow with invertible 1x1 convolutions. In *Advances in Neural Information Processing Systems*, pp. 10236–10245, 2018.
- Le, T. A., Baydin, A. G., Zinkov, R., and Wood, F. Using synthetic data to train neural networks is model-based reasoning. In *Neural Networks (IJCNN), 2017 International Joint Conference on*, pp. 3514–3521. IEEE, 2017.
- Lotka, A. J. Analytical note on certain rhythmic relations in organic systems. *Proceedings of the National Academy of Sciences*, 6(7):410–415, 1920.
- Louppe, G. and Cranmer, K. Adversarial variational optimization of non-differentiable simulators. *arXiv preprint arXiv:1707.07113*, 2017.
- Lueckmann, J.-M., Goncalves, P. J., Bassetto, G., Öcal, K., Nonnenmacher, M., and Macke, J. H. Flexible statistical inference for mechanistic models of neural dynamics. In *Advances in Neural Information Processing Systems*, pp. 1289–1299, 2017.
- Lueckmann, J.-M., Bassetto, G., Karaletsos, T., and Macke, J. H. Likelihood-free inference with emulator networks. In Ruiz, F., Zhang, C., Liang, D., and Bui, T. (eds.), *Proceedings of The 1st Symposium on Advances in Approximate Bayesian Inference*, volume 96 of *Proceedings of Machine Learning Research*, pp. 32–53. PMLR, 2018.
- Marjoram, P., Molitor, J., Plagnol, V., and Tavaré, S. Markov chain monte carlo without likelihoods. *Proc Natl Acad Sci U S A*, 100(26), 2003.
- May, R. M. and Leonard, W. J. Nonlinear aspects of competition between three species. *SIAM journal on applied mathematics*, 29(2):243–253, 1975.
- McCarthy, A., Rodriguez, B., and Minchale, A. Variational inference over non-differentiable cardiac simulators using bayesian optimization. *arXiv preprint arXiv:1712.03353*, 2017.
- Mohamed, S. and Lakshminarayanan, B. Learning in implicit generative models. *ArXiv e-prints*, 2016.
- Ong, V. M.-H., Nott, D. J., Tran, M.-N., Sisson, S. A., and Drovandi, C. C. Likelihood-free inference in high dimensions with synthetic likelihood. *Computational Statistics & Data Analysis*, 128:271–291, 2018.
- Owen, J., Wilkinson, D. J., and Gillespie, C. S. Likelihood free inference for markov processes: a comparison. *Statistical applications in genetics and molecular biology*, 14(2):189–209, 2015.
- Paninski, L. and Cunningham, J. Neural data science: accelerating the experiment-analysis-theory cycle in large-scale neuroscience. *bioRxiv*, pp. 196949, 2017.
- Papamakarios, G. and Murray, I. Fast ε -free inference of simulation models with bayesian conditional density estimation. In *Advances in Neural Information Processing Systems*, pp. 1028–1036, 2016.
- Papamakarios, G., Pavlakou, T., and Murray, I. Masked autoregressive flow for density estimation. In *Advances in Neural Information Processing Systems*, pp. 2338–2347, 2017.
- Papamakarios, G., Sterratt, D. C., and Murray, I. Sequential neural likelihood: Fast likelihood-free inference with autoregressive flows. *arXiv preprint arXiv:1805.07226*, 2018.
- Pesah, A., Wehenkel, A., and Louppe, G. Recurrent machines for likelihood-free inference. *arXiv preprint arXiv:1811.12932*, 2018.
- Pritchard, J. K., Seielstad, M. T., Perez-Lezaun, A., and Feldman, M. W. Population growth of human y chromosomes: a study of y chromosome microsatellites. *Mol Biol Evol*, 16(12), 1999.
- Reichenbach, T., Mobilia, M., and Frey, E. Mobility promotes and jeopardizes biodiversity in rock–paper–scissors games. *Nature*, 448(7157):1046, 2007.
- Reichenbach, T., Mobilia, M., and Frey, E. Self-organization of mobile populations in cyclic competition. *Journal of Theoretical Biology*, 254(2):368–383, 2008.
- Rezende, D. J. and Mohamed, S. Variational inference with normalizing flows. *arXiv preprint arXiv:1505.05770*, 2015.
- Romaszko, L., Williams, C. K., Moreno, P., and Kohli, P. Vision-as-inverse-graphics: Obtaining a rich 3d explanation of a scene from a single image. In *ICCV workshops*, pp. 940–948, 2017.
- Sisson, S. A., Fan, Y., and Tanaka, M. M. Sequential monte carlo without likelihoods. *Proceedings of the National Academy of Sciences*, 104(6):1760–1765, 2007.
- Sisson, S. A., Fan, Y., and Beaumont, M. *Handbook of Approximate Bayesian Computation*. Chapman and Hall/CRC, 2018.
- Tran, D., Ranganath, R., and Blei, D. Hierarchical Implicit Models and Likelihood-Free Variational Inference. In *Advances in Neural Information Processing Systems 30*. Curran Associates, Inc., 2017a.
- Tran, M.-N., Nott, D. J., and Kohn, R. Variational bayes with intractable likelihood. *Journal of Computational and Graphical Statistics*, 26(4):873–882, 2017b.

Turner, B. M. and Sederberg, P. B. A generalized, likelihood-free method for posterior estimation. *Psychonomic Bulletin & Review*, 21(2), 2014.

Wood, S. N. Statistical inference for noisy nonlinear ecological dynamic systems. *Nature*, 466(7310):1102, 2010.

Zaheer, M., Kottur, S., Ravanbakhsh, S., Poczos, B., Salakhutdinov, R. R., and Smola, A. J. Deep sets. In *Advances in neural information processing systems*, pp. 3391–3401, 2017.

A. Supplementary material

A.1. Derivation of MoG proposals for APT with MDNs

When $q_{F(x,\phi)}$ is an M -component mixture of Gaussians, $\tilde{p}(\theta)$ is an L -component mixture of Gaussians and $p(\theta)$ is Gaussian, we have

$$q_{F(x,\phi)}(\theta) = \sum_{i=1}^M \alpha_i \mathcal{N}_\theta(\mu_i, \Sigma_i) \quad (7)$$

$$\tilde{p}(\theta) = \sum_{k=1}^L \beta_k \mathcal{N}_\theta(\tilde{\mu}_k, \tilde{\Sigma}_k) \quad (8)$$

$$p(\theta|x) = \mathcal{N}_\theta(\mu_0, \Sigma_0) \quad (9)$$

$$\tilde{q}_{x,\phi}(\theta) = \frac{1}{Z_{x,\phi}} \sum_{i,k} \alpha_i \beta_k \frac{\mathcal{N}_\theta(\mu_i, \Sigma_i) \mathcal{N}_\theta(\tilde{\mu}_k, \tilde{\Sigma}_k)}{\mathcal{N}_\theta(\mu_0, \Sigma_0)} \quad (10)$$

$$= \sum_{i,k} \zeta_{ik} \mathcal{N}_\theta(\mu_{ik}^*, \Sigma_{ik}^*) \quad (11)$$

where

$$\Sigma_{ik}^* = \left(\Sigma_i^{-1} + \tilde{\Sigma}_k^{-1} - \Sigma_0^{-1} \right)^{-1} \quad (12)$$

$$\mu_{ik}^* = \Sigma_{ik}^* \left(\Sigma_i^{-1} \mu_i + \tilde{\Sigma}_k^{-1} \tilde{\mu}_k - \Sigma_0^{-1} \mu_0 \right) \quad (13)$$

$$\zeta_{ik} \propto \alpha_i \beta_k \sqrt{\frac{\det(\Sigma_{ik}^*)}{\det(\Sigma_i) \det(\tilde{\Sigma}_k)}} e^{-\frac{1}{2} (\mu_i^\top \Sigma_i^{-1} \mu_i + \tilde{\mu}_k^\top \tilde{\Sigma}_k^{-1} \tilde{\mu}_k - \mu_{ik}^{*\top} \Sigma_{ik}^{-1} \mu_{ik}^*)} \quad (14)$$

and the proportionality symbol indicates that the weights ζ_{ik} should be normalized so that $\sum_{ik} \zeta_{ik} = 1$.

A.2. Algorithm and computational complexity for atomic APT

Algorithm 2 APT with atomic proposals

Input: simulator with (implicit) density $p(x|\theta)$, data x_o , prior $p(\theta)$, density family q_ψ , neural network $F(x, \phi)$, simulations per round N , number of rounds R , number of atoms M .

```

 $\tilde{p}_1(\theta) := p(\theta)$ 
 $c \leftarrow 0$  total simulation count
for  $r = 1$  to  $R$  do
    for  $j = 1$  to  $N$  do
         $c \leftarrow c + 1$ 
        Sample  $\theta_c \sim \tilde{p}_r(\theta)$ 
        Simulate  $x_c \sim p(x|\theta_c)$ 
    end for
     $V_r(\Theta) := \begin{cases} \binom{c}{M}^{-1}, & \text{if } \Theta = \{\theta_{b_1}, \theta_{b_2}, \dots, \theta_{b_M}\} \text{ and } 1 \leq b_1 < b_2 < \dots < b_M \leq c \\ 0, & \text{otherwise} \end{cases}$  sampling without replacement
     $\phi \leftarrow \operatorname{argmin}_\phi \mathbb{E}_{\Theta \sim V_r(\Theta)} \left[ \sum_{\theta_j \in \Theta} -\log \tilde{q}_{x_j, \phi}(\theta_j) \right]$ 
     $\tilde{p}_{r+1}(\theta) := q_{F(x_o, \phi)}(\theta)$ 
end for
return  $q_{F(x_o, \phi)}(\theta)$ 

```

To minimize $\mathbb{E}_{\Theta \sim V_i(\Theta)} \left[\sum_{\theta_j \in \Theta} -\log \tilde{q}_{x_j, \phi}(\theta_j) \right]$, we must calculate its gradient with respect to ϕ , for which we use mini-batches of size M . Specifically, for each minibatch we first sample $B = \{b_1, \dots, b_M\} \subset \{1, \dots, c\}$ without replacement. We then calculate the gradient $\frac{d}{d\phi} \sum_{b \in B} -\log \tilde{q}_{x_b, \phi}(\theta_b)$ using (6). Note that each minibatch involves M loss evaluations—that is, M multiple-choice questions with M possible answers each.

To calculate $\tilde{q}_{x_j, \phi}$ and its gradients for a single minibatch using (6), we must evaluate $q_{F(x, \phi)}(\theta)$ on every possible pair $(\theta_b, x_{b'})$ for $b, b' \in B$. Therefore atomic APT's computational complexity is quadratic in the minibatch size M . However, we observed no difference in wallclock time per minibatch for $M = 10$ vs. $M = 100$ on an nVidia GeForce RTX 2080. Furthermore, for the Lokta-Volterra and RPS models, simulations took longer than all other calculations, for all methods.

When using an MDN as the density estimator, the MoG estimate of the posterior can be calculated once for each $x_{b'}$, and then each MoG evaluated on each θ_b . For a network with n_{layers} fully connected hidden layers of n_{hidden} units each, $d = \dim(x)$ and an n_{MoG} -component Gaussian mixture, the computational complexity of an atomic APT minibatch is

$$C_{\text{atomic MDN-APT}} = \mathcal{O}(M [dn_{\text{hidden}} + n_{\text{layers}}n_{\text{hidden}}^2 + n_{\text{hidden}}n_{\text{MoG}}\dim(\theta)^2] + M^2n_{\text{MoG}}\dim(\theta)^2) \quad (15)$$

Note that only the final term is quadratic in M , and this term does not involve the network structure or input dimensionality. For comparison, SNPE-A/B or non-atomic MDN-based APT has the same complexity, except for being linear in M :

$$C_{\text{SNPE-A/B/non-atomic MDN-APT}} = \mathcal{O}(M [dn_{\text{hidden}} + n_{\text{layers}}n_{\text{hidden}}^2 + n_{\text{hidden}}n_{\text{MoG}}\dim(\theta)^2]) \quad (16)$$

In a MAF data and parameters are coupled, so for atomic APT every $(\theta_b, x_{b'})$ pair requires a separate feedforward pass. For n_{MADEs} conditional MADEs consisting of n_{layers} fully connected hidden layers of n_{hidden} units, the minibatch complexity is

$$C_{\text{atomic MAF-APT}} = \mathcal{O}(n_{\text{MADEs}} [Md n_{\text{hidden}} + M^2 \dim(\theta)n_{\text{layers}}n_{\text{hidden}}^2 + M^2 \dim(\theta)^2 n_{\text{hidden}}]) \quad (17)$$

The first term is exempted from quadratic dependence on M since the inputs x can be multiplied by the appropriate weight matrices independently of θ . Thus, any computational tasks that scale with the input dimension d scale only linearly with M .

For MAF-based SNL the roles of the data and parameters are reversed, so the complexity (not including MCMC sampling) is

$$C_{\text{MAF-SNL}} = \mathcal{O}(Mn_{\text{MADEs}} [\dim(\theta)n_{\text{hidden}} + dn_{\text{layers}}n_{\text{hidden}}^2 + d^2n_{\text{hidden}}]) \quad (18)$$

For very large d , we hypothesize that atomic MAF-APT could benefit from specialized network architectures that pass x through a learned, feed-forward dimensionality reduction before supplying the result as input to all MADEs.

A.3. Conditional flow normalization and truncated priors

Conditional density estimators trained to target the posterior density should respect constraints imposed on them by the prior. For uniform priors, the posterior density outside the prior support should be zero. The proposal correction for SNPE-A (Papamakarios & Murray, 2016) does not hold for tight priors that clearly truncate the posterior estimate. Also for SNPE-B, posterior leakage leads to hard-to-interpret results unless the MDNs are re-normalized.

The normalization of the conditional density model $q_{F(x,\phi)}(\theta)$ cancels out in the computation of the probabilities $\tilde{q}_{x,\phi}(\theta)$ (see eq. 6). Hence conditional density models optimized via APT with atomic proposals are automatically normalized during training. After training, the Gaussian (mixture) $q_{F(x_o,\phi)}(\theta)$ returned from MDNs can be truncated to return a valid truncated Normal posterior estimate. For MAFs we can effectively obtain such post-hoc truncation through rejection sampling, and estimate the normalization factor from the rejection rates.

We noticed that across many rounds, conditional MAFs trained with APT can leak increasingly large amounts of mass outside the prior support. We did not find this leakage to negatively influence the quality of the estimated posterior shape (evaluated over the prior support). If needed, there would be several options to reduce leakage across rounds and keep rejection rates low: We can periodically reinitialize the conditional density estimator across rounds. Alternatively, one could also train a new flow which is optimized to have the same shape on the prior’s support, but minimal mass elsewhere. This normalized flow could then be used to directly evaluate posterior densities. For simple box-shaped prior supports, one can also apply a pointwise (scaled) logistic transformation to the MAF outputs to enforce prior bounds, i.e. train $q_{F(x,\phi)}(\sigma^{-1}(\theta))$.

A.4. Proof of proposition 1

Here we prove Prop. 1. Note that whenever we refer to $\tilde{p}(\theta)$, $\tilde{p}(\theta|x)$ or $\tilde{p}(x)$ these distributions are always defined based on some specific choice of Θ .

Proof of proposition 1

$$\mathbb{E}_{\Theta \sim V, \theta \sim U_\Theta, x \sim p(x|\theta)} [\tilde{\mathcal{L}}(\phi) - \tilde{\mathcal{L}}(\phi^*)] = \int_{\Theta} V(\Theta) \sum_{\theta \in \Theta} \tilde{p}(\theta) \int_x p(x|\theta) [\log \tilde{p}(\theta|x) - \log \tilde{q}_{x,\phi}(\theta)] \quad (19)$$

By Bayes’ rule and (1), for $\theta \in \Theta$ we have $p(x|\theta) = \frac{\tilde{p}(\theta|x)\tilde{p}(x)}{\tilde{p}(\theta)}$ so

$$\mathbb{E}_{\Theta \sim V, \theta \sim U_\Theta, x \sim p(x|\theta)} [\tilde{\mathcal{L}}(\phi) - \tilde{\mathcal{L}}(\phi^*)] = \int_{\Theta} V(\Theta) \int_x \tilde{p}(x) \sum_{\theta \in \Theta} \tilde{p}(\theta|x) \log \frac{\tilde{p}(\theta|x)}{\tilde{q}_{x,\phi}(\theta)} \quad (20)$$

$$= \int_{\Theta} V(\Theta) \int_x \tilde{p}(x) D_{\text{KL}}(\tilde{p}(\theta|x) || \tilde{q}_{x,\phi}(\theta)) \quad (21)$$

$$= \mathbb{E}_{\Theta \sim V, \theta \sim U_\Theta, x \sim p(x|\theta)} [D_{\text{KL}}(\tilde{p}(\theta|x) || \tilde{q}_{x,\phi}(\theta))] \quad (22)$$

By Gibbs’ inequality, the KL divergence is zero only when $\tilde{q}_{x,\phi}(\theta) = \tilde{p}(\theta|x)$ for all $\theta \in \Theta$, in which case by (5-6) $q_{F(x,\phi)} \propto p(\theta|x)$ for $\theta \in \Theta$ as well. Thus $D_{\text{KL}}(\tilde{p}(\theta|x) || \tilde{q}_{x,\phi}(\theta)) > 0$ whenever $\rho(x, \Theta, \phi) \neq 1$, so if $\mathbb{E}_{\Theta \sim V, \theta \sim U_\Theta, x \sim p(x|\theta)} [\rho(x, \Theta, \phi)]$ were less than one, $\mathbb{E}_{\Theta \sim V, \theta \sim U_\Theta, x \sim p(x|\theta)} [\tilde{\mathcal{L}}(\phi) - \tilde{\mathcal{L}}(\phi^*)]$ would be greater than zero.

A.5. Experimental details

We use the same basic network architectures for all of our experiments. For mixture-density networks (SNPE-A, SNPE-B, APT), we use two fully-connected tanh layers with 50 units each. Unless otherwise stated, we use MDNs with 8 Gaussian mixture components. In our experiments with MDNs MoG proposals were used for APT. For conditional masked autoregressive flows (SNL, APT), we use stacks of 5 MADEs each constructed using two fully-connected tanh layers with 50 units each. We train the APT MAFs with atomic proposals using $M = 100$ atoms.

For the SLCP, Lotka-Volterra and M/G/1 model, we follow the experimental setup of (Papamakarios et al., 2018), including uniform priors, summary statistics, ground-truth parameters θ^* and observed data x_o .

A.5.1. TWO MOONS MODEL

For given parameter $\theta \in \mathbb{R}^2$, the ‘two moons’ simulator generates $x \in \mathbb{R}^2$ according to

$$a \sim U\left(-\frac{\pi}{2}, \frac{\pi}{2}\right) \quad (23)$$

$$r \sim \mathcal{N}(0.1, 0.01^2) \quad (24)$$

$$p = (r \cos(a) + 0.25, r \sin(a)) \quad (25)$$

$$x^\top = p + \left(-\frac{|\theta_1 + \theta_2|}{\sqrt{2}}, \frac{-\theta_1 + \theta_2}{\sqrt{2}} \right) \quad (26)$$

The intermediate variables p follow a single crescent-shaped distribution, which is subsequently shifted and rotated around the origin depending on θ . Consequently, $p(x|\theta)$ shows a single shifted crescent for fixed θ . The absolute value $|\theta_1 + \theta_2|$ gives rise to the second crescent in the posterior. We choose a uniform prior over $[-1, 1]^2$ to illustrate inference on this model. As observed data we use $x_o = (0, 0)^\top$.

On this example we fit mixture-density networks with 20 mixture components to allow expressive conditional densities.

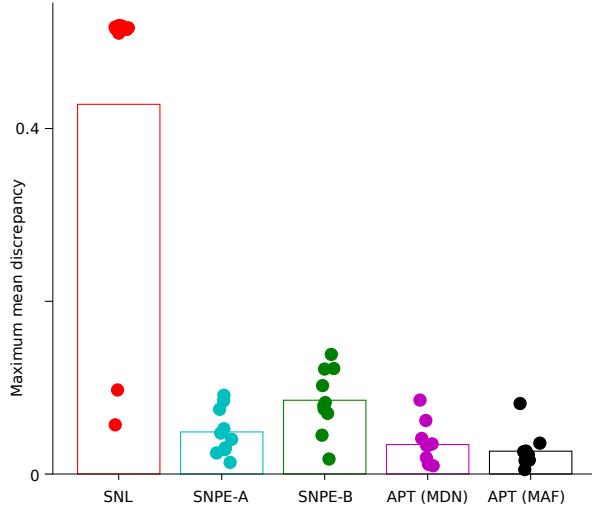


Figure 7. Two moons model. Comparison of average maximum mean discrepancies between final posterior estimate and ground-truth posteriors for different algorithms across 10 different random initializations and $x = (0, 0)$. Dots show individual runs. The MCMC chains used by SNL to obtain posterior estimates failed to sample both posterior modes in the majority of cases (cf. Fig 1 for such an example case), leading to high discrepancies.

A.5.2. SLCP MODEL

The SLCP model has a simple simulator density $p(x|\theta) = \prod_{i=1}^4 \mathcal{N}(x_{(2i-1, 2i)} | \mu(\theta), \Sigma(\theta))$, i.e. $x \in \mathbb{R}^8$ consists of four independent samples from a bivariate Gaussian parameterized by $\theta \in \mathbb{R}^5$. The conditional mean is given by $\mu(\theta) = (\theta_1, \theta_2)^\top$. The parameterization of the covariance

$$\Sigma(\theta) = \begin{bmatrix} \theta_3^2 & \theta_3 \theta_4 \tanh(\theta_5) \\ \theta_3 \theta_4 \tanh(\theta_5) & \theta_4^2 \end{bmatrix} \quad (27)$$

leads to in total four modes in $p(\theta|x)$ visible in the pairwise marginal over (θ_3, θ_4) (Fig. 2a). To calculate MMD’s, We sampled from the ground-truth posterior using MCMC. For the experiment with added uninformative simulator outputs, we generate noise outputs $x_{i>8} \in \mathbb{R}^m$, $m \in \{12, 32, 52, 92\}$ from m -dimensional mixtures of t-distributions and append them to the 8-dimensional simulator output. We use mixtures of 20 multivariate t-distributions with randomized means and covariance matrices and degree of freedom 2 to create non-trivial densities $p(x|\theta)$. We note that the actual posterior density $p(\theta|x)$ (for any noise outputs $x_{i>8}$) retains the shape of the original SLCP model due to

$$p(\theta|x) = \frac{p(x|\theta)p(\theta)}{p(x)} = \frac{p(x_{i>8}|x_{i\leq 8}, \theta)p(x_{i\leq 8}|\theta)p(\theta)}{p(x_{i>8}|x_{i\leq 8})p(x_{i\leq 8})} = \frac{p(x_{i>8})p(x_{i\leq 8}|\theta)p(\theta)}{p(x_{i>8})p(x_{i\leq 8})} = \frac{p(x_{i\leq 8}|\theta)p(\theta)}{p(x_{i\leq 8})} \quad (28)$$

To avoid effects of the autoregressive nature of MAF density estimators, we re-order the dimensions i of the original eight output dimensions $x_{i\leq 8}$ and our added simulator outputs $x_{i>8}$ with a fixed random permutation.

Note that in Fig. 3 the MMD for the true posterior (lower gray lines) is nonzero due to a finite number of samples being used.

A.5.3. LOTKA-VOLTERRA MODEL

For the comparisons against previous neural conditional density models, we apply APT to infer the posterior of the Lotka-Volterra model as described in (Papamakarios et al., 2018).

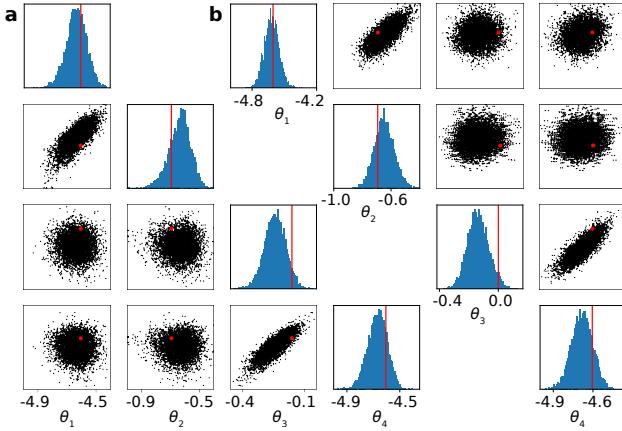


Figure 8. Lotka-Volterra model. Close-up comparison between **a** APT and **b** SNL posterior estimates.

We also infer posteriors for a Lotka-Volterra model with added observation noise. We add independent Gaussian noise $\varepsilon_{it} \sim \mathcal{N}(0, \sigma^2)$ to both populations $i = 1, 2$ and every time point $t = 0, \dots, 150$ of the raw simulated time-series. For RNN-APT, we added an initial layer of 100 GRU units (Cho et al., 2014) to a MDN with a single Gaussian component.

A.5.4. ROCK-PAPER-SCISSORS MODEL

We approximated the SPDE using a system of coupled stochastic differential equations as described in eq. (29-30) of (Reichenbach et al., 2008), and integrated this system using the Euler-Murayama method with a step size of 1 on a 100x100 grid. We calculated second spatial derivatives using simple finite differences-of-differences. We initialized the system at the unstable uniform steady state, and integrated from $t = 0$ to $t = 100$. With μ , σ and D denoting the growth rate, predation rate and diffusion constant as defined in (Reichenbach et al., 2007; 2008), θ_1 , θ_2 and θ_3 were defined as their respective base 10 logarithms. We used uniform priors on each θ_j , from -1 to 1 for $1 \leq j \leq 2$ and from -6 to -5 for $j = 3$.

We used a CNN consisting of 6 convolutional layers with ReLu units, with each convolutional layer followed by a max pooling layer. The number of channels after each layer was: 8, 8, 8, 16, 32 and 32. The convolutional filter sizes were: 3, 3, 3, 3, 2 and 2. The max pooling sizes were: 1, 3, 2, 2, 2 and 1. The image sizes were 100x100, 32x32, 15x15, 6x6, 2x2 and 1x1. The 32-dimensional output of the CNN was then passed through two fully-connected tanh layers with 50 units each. We used a single-component MDN for this problem. Overall, for this architecture the elements of ϕ (weights and biases) numbered 8768 for the CNN layers, 4200 for the fully connected layers and 612 for the final mapping onto ψ .

A.5.5. M/G/1 MODEL

We compared MAF-APT to previous approaches on the M/G/1 queue model (as described in Papamakarios & Murray, 2016). The results (mean \pm SEM) for 10 different random initializations with the identical x_o are shown in figure 9.

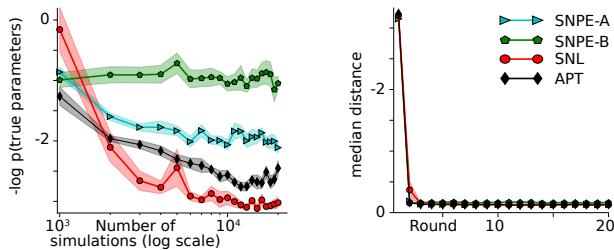


Figure 9. M/G/1 model. Averages ± 1 SEM over 10 different random intializations with identical x_o (≤ 10 for SNPE-A in later rounds).

A.5.6. GENERALIZED LINEAR MODEL

We also compare MAF-APT against previous neural conditional density approaches on a Generalized Linear model model with a length-9 temporal input filter and a bias weight (i.e. $\theta_j, j = 1, \dots, d$ with $d = 10$ parameters in total). We simulate 100 time bins of activity in response to white noise and summarize the output with 10 sufficient statistics x as in (Lueckmann et al., 2017). We also train the algorithms with 5 rounds of $N = 5000$ each. We note that the posteriors for this simulator are well approximated as Gaussian, and use a single Gaussian component for the MDNs for SNPE-A and -B, and MAFs consisting of two MADES for SNL and APT. We use networks with two layers of 10 tanh units for MDNs and MADES.

The results (mean \pm SEM) for 10 different random initializations and x_o are shown in figure 10.

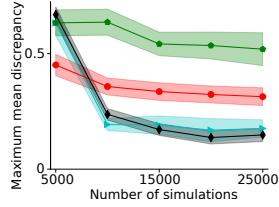


Figure 10. Generalized linear model. Averages \pm 1 SEM over 10 different random intializations.