

Amortizing Bayesian Posterior Inference in Tractable Likelihood Models

Anonymous authors
Paper under double-blind review

Abstract

Bayesian inference provides a natural way of incorporating prior beliefs and assigning a probability measure to the space of hypotheses. However, it is often infeasible in practice as it requires expensive iterative routines like MCMC to approximate the posterior distribution. Not only are these methods computationally expensive, but they must also be re-run whenever new observations are available, making them impractical or of limited use. To alleviate such difficulties, we amortize the posterior parameter inference for probabilistic models through permutation invariant architectures. While this paradigm is briefly explored in Simulation Based Inference (SBI), Neural Processes (NPs) and Gaussian Process (GP) kernel estimation, a more general treatment of amortized Bayesian inference in known likelihood models has been largely unexplored. We rely on the reverse-KL based amortized Variational Inference (VI) approach to train inference systems and compare them with forward-KL based SBI approaches across different architectural setups. We additionally introduce a simple but strong approach to further amortize on the number of features in each observation, allowing a single system to infer variable dimensional parameters. Our thorough experiments demonstrate the effectiveness of our proposed approach, especially in real-world and model misspecification settings.

1 Introduction

Bayesian analysis of data has become increasingly popular and is widely used in numerous scientific disciplines. In politics, predictive models based on public polling and other factors play a crucial role in the discourse around the state of a campaign. Throughout the COVID-19 pandemic, models that estimate the infectiousness of the virus, the efficacy of public health measures, and the future course of the pandemic became critical to government planning and the public’s understanding of the pandemic (Cooper et al., 2020). In cryogenic electron microscopy (cryo-EM), the posterior over an unknown 3D atomic-resolution molecular structure is explored given the 2D image observations (Glaeser et al., 2021).

While recent years have made such methods more accessible (Bingham et al., 2019; Carpenter et al., 2017; Štrumbelj et al., 2023), they still remain computationally burdensome. Further, in practical contexts where new observations are continuously available, the analysis must be re-run every time new data becomes available, e.g., when new case counts become available, previous measurements are corrected, or when applied to different geographic regions. As a result practitioners adopt approximations (Welling & Teh, 2011; Gelfand, 2000; Brooks, 1998), simplify their models (Hoffman et al., 2013; Blei et al., 2017) or reduce the frequency with which they perform their analyses.

A common thread is that the probabilistic model defining the relationship between its parameters and the observations is fixed. Poll aggregation models use hierarchical time series models (Athanasopoulos et al., 2023; Chen et al., 2023), infectious diseases are studied using variants on compartment models (Tang et al., 2020), and cryo-EM uses a linear image formation model (Glaeser et al., 2021). This makes these applications ideal candidates for amortized inference (Morris, 2013; Paige & Wood, 2016; Kingma & Welling, 2013; Rezende et al., 2014; Stuhlmüller et al., 2013).

We propose using neural networks to learn a function that maps an observed *dataset* directly to the corresponding posterior distribution without the need for iterative procedures, e.g., Markov chain Monte

Carlo (MCMC) sampling (Gelfand, 2000; Hoffman et al., 2014). To efficiently handle permutation invariance stemming from the *iid* nature of observations, we rely on efficient set-based architectures like Transformers and DeepSets (Zaheer et al., 2017; Vaswani et al., 2017; Lee et al., 2019). If learned properly, this mapping would allow generalization to different datasets zero-shot. We also leverage a simple masking procedure to amortize posterior estimation over datasets with a variable number of features, thus inching closer towards a generalist Bayesian learner – as evidenced by its generalization capabilities on real-world tasks zero-shot through only pre-training on synthetic data.

Generally, real-world datasets do not exactly follow standard models, e.g., while practitioners often rely on linear models, data rarely follows them exactly. Amortization methods like simulation-based inference (SBI), in particular neural posterior estimation (SBI-NPE), (Radev et al., 2020; Geffner et al., 2023; Cranmer et al., 2020) suffer from a gap in the training (simulation) and evaluation (real-world) protocol and thus can be suboptimal. Even further, they require access to the underlying model parameters generating the datasets during training and are thus not applicable in many practical settings where the true parameters are unknown.

Our primary motivation is posterior inference as the parametric values themselves are often of interest in applied statistical practice, e.g., for assessing the success of a pandemic intervention or the impact of a factor on public opinion polling. However, evaluating approximate posteriors is extremely difficult since the true posterior is highly intractable to either evaluate or sample from with the exception of MCMC sampling which only provides asymptotic guarantees. Despite their intractability, Bayesian posteriors have been claimed to provide better generalization capabilities as they encompass multiple different hypotheses in making predictions. Following this line of reasoning, we demonstrate the utility of our proposed approach in the closely related problem of posterior prediction, where the goal is to model future predictions given some observations.

Since SBI-NPE approaches are generally modeled through minimization of forward-KL objective, they suffer from over-dispersion during Bayesian posterior inference and consequently worse predictive performance. Instead, we propose a new training objective that can operate solely on datasets without access to the underlying parameters, allowing for a wider diversity of data to be incorporated during training for better generalization to real world settings. Our detailed experiments and benchmarking establish the superiority of our proposed approach as well as their qualitative differences. Our contributions include

- Proposing a general framework for amortizing Bayesian posterior estimation in probabilistic models and highlighting its superior performance to SBI-NPE.
- Extending amortization in both SBI-NPE and the proposed framework over a variable number of features in each observation, and not just the number of observations.
- Quantifying benefits of our proposed approach when the true underlying model is unknown (model misspecification), as well as providing its advantages in generalization to real-world tabular datasets.
- Benchmarking various design choices like architectural backbones through extensive ablations.

2 Background

We first cover some of the important preliminaries as well as approaches already existing in the literature.

Bayesian Inference. Let $\mathbf{x} \in \mathbb{R}^d$ denote the outcome of an experiment observed through a set of independent and identically (*iid*) distributed samples $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$. Given these observations, we are often interested in either quantifying the certainty of or generating potential future observations \mathbf{x}_* . Bayesian Inference provides a natural methodology of quantifying $p(\mathbf{x}_*|\mathcal{D})$ by prescribing a space of hypotheses $\boldsymbol{\theta} \in \mathbb{R}^k$ and a *prior* belief $p(\boldsymbol{\theta})$ over it. These hypotheses $\boldsymbol{\theta}$ define the *likelihood* of observing a particular outcome, i.e., $p(\mathbf{x}|\boldsymbol{\theta})$. The likelihood and prior are then combined through Bayes rule to define the *posterior* $p(\boldsymbol{\theta}|\mathcal{D})$, through which the quantity of interest can then be easily expressed as

$$p(\mathbf{x}_*|\mathcal{D}) = \int_{\boldsymbol{\theta}} p(\mathbf{x}_*|\boldsymbol{\theta})p(\boldsymbol{\theta}|\mathcal{D})d\boldsymbol{\theta} \tag{1}$$

However, the above expression poses two challenges: (a) the *posterior* $p(\boldsymbol{\theta}|\mathcal{D})$, which is often a quantity of interest in itself, is not known, and (b) even if known, the integration might be intractable. The intractability

of the integration is often resolved through Monte Carlo estimation

$$p(\mathbf{x}_*|\mathcal{D}) = \mathbb{E}_{\boldsymbol{\theta}|\mathcal{D}} [p(\mathbf{x}_*|\boldsymbol{\theta})] \quad (2)$$

$$\approx \frac{1}{M} \sum_{m=1}^M p(\mathbf{x}_*|\boldsymbol{\theta}^{(m)}) \quad (3)$$

where $\boldsymbol{\theta}^{(m)} \sim p(\boldsymbol{\theta}|\mathcal{D})$. The quantity $p(\boldsymbol{\theta}|\mathcal{D})$ can be obtained through an application of Bayes rule

$$p(\boldsymbol{\theta}|\mathcal{D}) = \frac{p(\mathcal{D}|\boldsymbol{\theta}) p(\boldsymbol{\theta})}{p(\mathcal{D})} \quad (4)$$

$$= \frac{p(\boldsymbol{\theta})}{p(\mathcal{D})} \prod_{n=1}^N p(\mathbf{x}_n|\boldsymbol{\theta}) \quad (5)$$

Given the form of the *likelihood* and the *prior*, the above distribution is often difficult to sample from, especially with the added complexity of the marginal $p(\mathcal{D}) = \int_{\boldsymbol{\theta}} p(\mathcal{D}|\boldsymbol{\theta}) p(\boldsymbol{\theta})$ being intractable. Additionally, the posterior itself is often of interest on its own, especially in cases where $\boldsymbol{\theta}$ is interpretable, for example, if we model the bias of a coin based on multiple tosses. We refer the readers to [Bishop & Nasrabadi \(2006\)](#) for applications of Bayesian Inference to supervised learning, etc.

Variational Inference. To bypass the intractability of the posterior distribution, or at least the difficulty to sample from it, VI methods approximate the true posterior $p(\boldsymbol{\theta}|\mathcal{D})$ with a variational distribution $q_{\varphi}(\boldsymbol{\theta})$ and convert the estimation problem into the following optimization problem

$$\varphi^* = \arg \min_{\varphi} \mathbb{KL}[q_{\varphi}(\cdot)||p(\cdot|\mathcal{D})] \quad (6)$$

which is equivalent to optimizing the Evidence Lower-Bound (ELBO) ([Gelman et al., 2013](#))

$$\varphi^* = \arg \max_{\varphi} \mathbb{E}_{\boldsymbol{\theta} \sim q_{\varphi}(\cdot)} \left[\log \frac{p(\mathcal{D}, \boldsymbol{\theta})}{q_{\varphi}(\boldsymbol{\theta})} \right] \quad (7)$$

The above optimization procedure finds a member in the family of variational distributions $\{q_{\varphi}\}_{\varphi}$ that is closest to the true posterior under the *reverse-KL* divergence. Once the optimal parameters φ^* are obtained, the *posterior predictive* distribution $p(\mathbf{x}_*|\mathcal{D})$ can be approximated as

$$p(\mathbf{x}_*|\mathcal{D}) \approx \mathbb{E}_{q_{\varphi^*}(\boldsymbol{\theta})} [p(\mathbf{x}_*|\boldsymbol{\theta})] \quad (8)$$

The family of distributions q_{φ} is chosen such that it is easy to sample from. Typical choices include independent multivariate Gaussian distribution (mean-field approximation) or normalizing flows ([Rezende & Mohamed, 2015](#); [Papamakarios et al., 2021](#); [Ardizzone et al., 2018-2022](#)).

Amortization. One of the most powerful capabilities of neural networks is their ability to learn and generalize to a wide variety of domains and settings provided sufficient variability during training. For example, Variational Autoencoders (VAEs) define a latent-variable model $p(\mathbf{x}, \mathbf{z})$ where \mathbf{x} represents the observation and \mathbf{z} the latent variable. VI typically relies on solving a separate optimization problem $q_{\varphi_i^*}(\mathbf{z}_i)$ for each posterior $p(\mathbf{z}_i|\mathbf{x}_i)$. The cost of learning separate variational approximations can be amortized through training of a joint network $q_{\varphi}(\mathbf{z}|\mathbf{x})$, where φ now represents the parameters of a neural network which takes \mathbf{x} explicitly as input. The VI procedure then reduces to optimizing φ , which is shared across all observations, as opposed to optimizing for separate φ_i 's, in the hope that $q_{\varphi}(\mathbf{z}_i|\mathbf{x}_i) \approx q_{\varphi_i}(\mathbf{z}_i)$ for any \mathbf{x}_i . When modeling using Gaussian distributions, this distinction can be seen as $q_{\varphi}(\mathbf{z}_i|\mathbf{x}_i) := \mathcal{N}(\cdot; \boldsymbol{\mu}_{\varphi}(\mathbf{x}_i), \boldsymbol{\Sigma}_{\varphi}(\mathbf{x}_i))$ while $q_{\varphi_i}(\mathbf{z}_i) := \mathcal{N}(\cdot; \boldsymbol{\mu}_{\varphi_i}, \boldsymbol{\Sigma}_{\varphi_i})$ (note the functional dependencies). In a similar fashion, [Garnelo et al. \(2018b\)](#) amortize on datasets as explicit inputs, while score-based generative models ([Song et al., 2020](#)) amortize on timesteps. Such models are largely successful owing to the generalization capabilities of neural networks to new unseen observations as long as the encoder $q_{\varphi}(\mathbf{z}|\mathbf{x}_i)$ is trained on diverse enough observations \mathbf{x}_i 's.

Simulation-Based Inference. SBI considers the problem of inferring the parameters of the simulator from observations. This is often tackled via neural posterior estimation methods (SBI-NPE) where a deep-learning

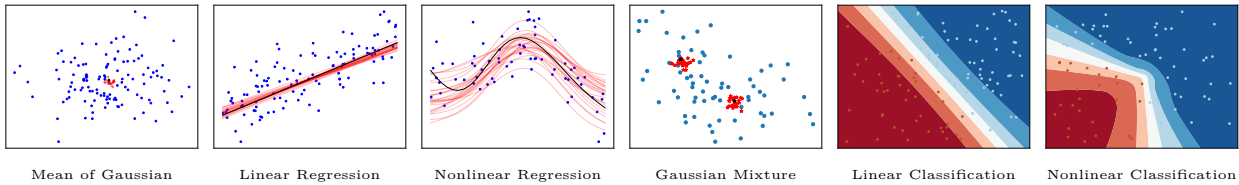


Figure 1: **Amortized Bayesian Posterior Estimation**: Illustration of predictions from proposed approach when trained on a fixed dimensional observation space. Model predictions, true predictions and sample points are shown in red, black and blue respectively. Additionally for classification, we label sample points with their ground-truth class, and draw the decision boundary according to the model.

based model is trained to infer the posterior by explicitly conditioning an approximate distribution $q_\varphi(\boldsymbol{\theta}|\mathcal{D})$ on the dataset, and modeling the gap between the distributions through a *Forward-KL* based optimization

$$\arg \min_{\varphi} \mathbb{E}_{\mathcal{D}} \text{KL} [p(\boldsymbol{\theta}|\mathcal{D}) || q_\varphi(\boldsymbol{\theta}|\mathcal{D})] \quad (9)$$

which often leads to mode averaging behavior that can be problematic in high dimensions. While the above objective often enjoys applications to tasks without tractable likelihood functions, it can only be used for training when the dataset \mathcal{D} is sampled according to the probabilistic model, and thus cannot utilize off-policy non-simulated data, hindering generalization to real-world scenarios. Precisely because of this, SBI-NPE has been leveraged in controlled scenarios like modeling inverse problems with low-dimensional non-differentiable simulators where the likelihood of an observation is not tractable; but has seen limited applicability in more general high dimensional estimation problems like the distribution over weights of a Bayesian Neural Networks.

Neural Processes. NPs also leverage amortized VI in training a latent-variable system for modeling predictive problems. However, unlike our setup, they define an approximate posterior distribution only over an arbitrary latent space and learn how this latent variable impacts the likelihood through point estimation of likelihood parameters. In particular, NPs rely on the Variational-EM setup, where they perform point estimation for the parameters of the likelihood and VI for the latent variable. In contrast, we focus on a similar setting where we instead do a full VI treatment of the parameters of the likelihood function; which in some sense means that our *latent variables* are now parameters of the likelihood model. Thus, our approach can be seen as a fully Bayesian Inference procedure for likelihood models, whereas NPs can be seen as latent-variable models which only provide point estimates for the parameters.

Gaussian Process Kernel Estimation. A specific application of our framework is the estimation of the kernel function for Gaussian Process likelihood models (Liu et al., 2020; Simpson et al., 2021; Bitzer et al., 2023), which leverages amortized inference for tractable likelihoods defined by GP regression setups. In contrast, we provide a more general framework for conducting amortized inference, which we test across a wide variety of domains ranging from supervised to unsupervised learning and from regression to classification tasks. Thus, our proposed approach provides a framework for parameter estimation through amortized variational inference and GP kernel estimation along these lines is a specific application of this approach.

In-Context Learning. ICL refers to large pre-trained models having the ability of solving new tasks at inference based solely on some context examples being provided as prompt. A natural parallel to conditioning on context examples via prompt in natural language is explicit amortization on observations for general-purpose tasks (Von Oswald et al., 2023; Müller et al., 2021). In this sense, our approach can be seen as training a general-purpose in-context learner to perform Bayesian inference, as opposed to direct predictions.

Sampling from Energy. Another line of similar works (Zhang & Chen, 2021; Berner et al., 2022; Vargas et al., 2023; Akhound-Sadegh et al., 2024; Bengio et al., 2021) look at sampling proportional to a pre-specified energy function, which in our case can be seen as $p(\boldsymbol{\theta}|\mathcal{D}) \propto e^{-\mathcal{E}(\boldsymbol{\theta})}$, with $\mathcal{E}(\boldsymbol{\theta}) = -\log p(\mathcal{D}, \boldsymbol{\theta})$. However, such works currently focus on non-amortized setups and we defer exploring amortized diffusion-styled posterior models as relevant future work.

We refer the readers to Appendix A for a detailed discussion about prior work, as well as its connections and differences with our proposed mechanism.

Objective q_φ	Model	L_2 Loss (\downarrow)					Accuracy (\uparrow)		
		Gaussian	GMM	LR	NLR		LC	NLC	
		$100D$	$5D$ $2 cl$	$100D$	$1D$	$25D$	$100D$	$2D$	$25D$
Baseline	- Random	301.06 \pm 0.35	5.00 \pm 0.04	202.6 \pm 0.3	65.94 \pm 0.91	831.6 \pm 8.7	50.0 \pm 0.0	50.3 \pm 0.6	50.0 \pm 0.3
	- Optimization	101.24 \pm 0.00	0.42 \pm 0.00	25.1 \pm 0.0	0.36 \pm 0.00	104.0 \pm 0.1	70.3 \pm 0.0	96.9 \pm 0.0	77.9 \pm 0.0
	- Langevin	102.35 \pm 0.03	0.45 \pm 0.01	23.3 \pm 0.7	0.31 \pm 0.00	132.4 \pm 1.0	65.1 \pm 0.4	96.0 \pm 0.3	73.2 \pm 0.3
	- HMC	102.41 \pm 0.03	0.48 \pm 0.01	18.7 \pm 0.2	0.37 \pm 0.00	98.1 \pm 0.7	62.1 \pm 0.2	91.8 \pm 0.2	70.4 \pm 0.1
Fwd-KL	GRU	102.64 \pm 0.01	2.43 \pm 0.03	124.8 \pm 0.1	49.33 \pm 0.95	671.6 \pm 10.5	59.7 \pm 0.1	59.5 \pm 0.4	56.9 \pm 0.3
	DeepSets	103.22 \pm 0.05	2.44 \pm 0.04	123.1 \pm 1.1	49.86 \pm 0.98	684.9 \pm 2.6	50.0 \pm 0.1	59.4 \pm 0.2	56.8 \pm 0.2
	Transformer	102.78 \pm 0.00	2.50 \pm 0.03	45.9 \pm 1.3	49.68 \pm 0.94	680.9 \pm 5.8	63.0 \pm 0.1	59.6 \pm 0.4	57.1 \pm 0.4
Fwd-KL	Gaussian GRU	102.51 \pm 0.01	0.47 \pm 0.01	60.2 \pm 0.9	0.43 \pm 0.00	106.0 \pm 0.6	63.5 \pm 0.3	92.4 \pm 0.2	72.5 \pm 0.0
	Gaussian DeepSets	102.60 \pm 0.04	0.50 \pm 0.02	62.8 \pm 0.6	0.43 \pm 0.00	125.9 \pm 0.8	60.9 \pm 0.3	92.5 \pm 0.1	59.8 \pm 0.3
	Gaussian Transformer	102.54 \pm 0.03	0.49 \pm 0.02	28.7 \pm 0.3	0.42 \pm 0.01	102.3 \pm 1.8	68.2 \pm 0.0	92.6 \pm 0.4	75.2 \pm 0.1
Fwd-KL	Flow GRU	102.66 \pm 0.02	0.67 \pm 0.09	119.1 \pm 0.2	15.78 \pm 0.21	539.0 \pm 4.3	59.9 \pm 0.2	76.9 \pm 0.3	58.3 \pm 0.0
	Flow DeepSets	103.34 \pm 0.03	0.65 \pm 0.08	125.7 \pm 3.7	15.05 \pm 0.12	548.5 \pm 3.3	50.1 \pm 0.0	72.3 \pm 1.8	58.1 \pm 0.1
	Flow Transformer	102.77 \pm 0.02	0.62 \pm 0.07	43.3 \pm 2.7	16.11 \pm 0.31	539.3 \pm 4.3	64.3 \pm 0.1	77.3 \pm 0.2	58.3 \pm 0.1
Rev-KL	GRU	102.49 \pm 0.01	0.47 \pm 0.00	61.3 \pm 1.0	0.41 \pm 0.01	106.0 \pm 0.4	64.7 \pm 0.2	93.4 \pm 0.1	72.0 \pm 0.5
	DeepSets	102.67 \pm 0.05	0.52 \pm 0.01	76.4 \pm 2.0	0.40 \pm 0.00	128.2 \pm 1.5	58.4 \pm 0.8	93.3 \pm 0.2	60.9 \pm 0.2
	Transformer	102.53 \pm 0.05	0.47 \pm 0.01	29.4 \pm 1.6	0.39 \pm 0.00	102.6 \pm 0.9	68.7 \pm 0.1	93.6 \pm 0.1	75.0 \pm 0.5

Table 1: **Amortized Bayesian Posterior Estimation.** Results for estimating the mean of a Gaussian (Gaussian), means of a Gaussian mixture model (GMM), (non-)linear regression (NLR/LR) and (non-)linear binary classification (NLC/LC). We ablate over different architectures and density parameterizations. Our baselines include the prior (Random), dataset-specific maximum likelihood training (Optimization), Langevin and Hamiltonian MCMC and Forward-KL based SBI. We use the expected L_2 loss and accuracy according to the posterior predictive as metrics.

3 Method

Given any modeling assumption, we are interested in estimating the posterior distribution over its parameters conditioned explicitly over observations to allow for a fast and scalable approximation. That is, given a probabilistic model $p(\cdot|\theta)$ with parameters θ , we are interested in approximating the full Bayesian inference solution over θ . Our goal is to train a system that approximates the posterior distribution $p(\theta|\mathcal{D})$ given a dataset $\mathcal{D} := \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\} \subseteq \mathbb{R}^d$ where $\mathbf{x}_i \sim p(\mathbf{x}|\theta)$. We achieve this by learning an amortized variational distribution $q_\varphi(\theta|\mathcal{D})$ conditioned explicitly on the dataset \mathcal{D} . Similar to standard VI approaches, we can train q_φ through the following optimization problem

$$\arg \min_{\varphi} \mathbb{KL} [q_\varphi(\cdot|\mathcal{D})||p(\cdot|\mathcal{D})] \quad (10)$$

While this is the case for VI on a single dataset, we are interested in generalizing to a family of datasets χ . Similar to how a VAE encoder efficiently generalizes to approximate posteriors for new images when trained across multiple image observations, we amortize over datasets by training over multiple $\mathcal{D} \sim \chi$ in the hope of generalizing to posterior approximations on new datasets. Thus, given a probability measure over the space of datasets χ , we posit learning of shared parameters φ of the variational distribution across different datasets as

$$\arg \min_{\varphi} \mathbb{E}_{\mathcal{D} \sim \chi} \mathbb{KL} [q_\varphi(\theta|\mathcal{D})||p(\theta|\mathcal{D})] \quad (11)$$

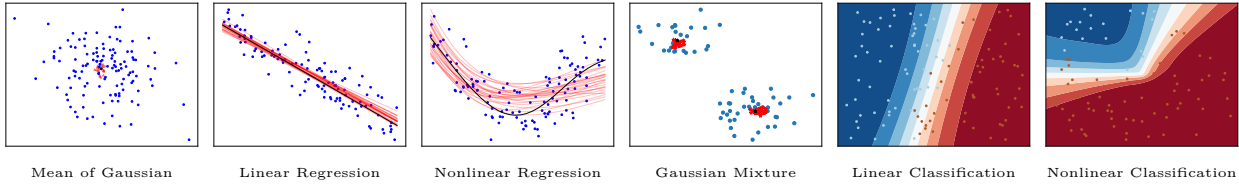


Figure 2: **Variable-Dimension Visualization:** Illustration of proposed approach when trained on a variable observational space but visualized on low-dimensional tasks, with model predictions, ground truth prediction and data points in red, black and blue respectively. For classification, data points are colored by their ground-truth labels, and decision boundary corresponds to the model prediction.

which equivalently reduces to maximizing the ELBO:

$$\arg \max_{\varphi} \mathbb{E}_{\mathcal{D} \sim \chi} \mathbb{E}_{\theta \sim q_{\varphi}(\cdot | \mathcal{D})} \left[\log \frac{p(\mathcal{D}, \theta)}{q_{\varphi}(\theta | \mathcal{D})} \right] \quad (12)$$

To recap, while VI strategies approach this problem by learning a separate q_{φ} for each \mathcal{D} , we instead train a shared q_{φ} that explicitly takes \mathcal{D} as input and, in doing so, possesses the ability to generalize to new datasets without additional training.

This training paradigm naturally introduces a dependency on the dataset generating distribution χ . Since we are working with a known probabilistic model, an obvious choice of χ is to treat this probabilistic model as a black-box simulator, akin to $p(n)p(\theta) \prod_{i=1}^n p(\mathbf{x}_i | \theta)$, samples from which can be obtained using ancestral sampling. Here, n is the dataset cardinality and $p(n)$ is a distribution over positive integers. Thus, obtaining at least one dataset-generating distribution is easy given any probabilistic model. However, χ can also be obtained from other sources, for example, a stream of real-world data, through interventions in the data-generating process, or through bagging on a large real-world dataset.

Another design choice in the above setup is the parameterization for q_{φ} : we explore a Gaussian distribution and a normalizing flow parameterization with either a Transformer, GRU or DeepSets architectural backbone to process \mathcal{D} . It is important to note that not all deep learning architectures (eg. GRU) are amenable in this setting since, given the *iid* nature of the samples, the posterior is permutation invariant

$$p(\theta | \mathcal{D}) = p(\theta | \Pi \mathcal{D}) \quad (13)$$

for some permutation matrix Π . This constraint should thus be reflected when modeling the approximation. To satisfy it, for example for a Gaussian distribution, we can define q_{φ} as

$$q_{\varphi}(\cdot) = \mathcal{N}(\cdot | \mu_{\varphi}(\mathcal{D}), \Sigma_{\varphi}(\mathcal{D})) \quad (14)$$

where μ_{φ} and Σ_{φ} are modeled using permutation invariant architectures, thus satisfying the desired constraint

$$q_{\varphi}(\cdot | \mathcal{D}) = q_{\varphi}(\cdot | \Pi \mathcal{D}) \quad (15)$$

Finally, another important design choice in this setup is the prior $p(\theta)$. In this work, we assume it to be known and fixed to $\mathcal{N}(\cdot; 0, \mathbf{I})$, which is a common assumption across various problems. Our proposed amortized VI approach, in contrast to SBI-NPE which relies on *Forward KL* optimization, relies on the *Reverse KL* objective.¹ Finally, we defer details regarding the architectural choices to Appendix B.

3.1 Amortizing Variable Feature Dimensions

So far, we have only considered amortization over different datasets for the same underlying likelihood model. For example, for a 2-dimensional Bayesian linear regression model, the amortized posterior $q_{\varphi}(\theta | \mathcal{D})$ approximates the true posterior distribution $p(\theta | \mathcal{D})$ for arbitrary sets of 2-dimensional observations. It is important to note that a deep learning-based approach leaves hopes of generalizing to new datasets since the underlying functional form of the solution remains constant across different datasets, irrespective of the number of features, and is given by the solution obtained from Equation 4.

¹Given the equivalences, we will use amortized VI and Reverse KL interchangeably.

Objective q_φ	Model	L_2 Loss (\downarrow)					Accuracy (\uparrow)		
		Gaussian	GMM	LR	NLR		LC	NLC	
		$100D$	$5D$ $2\ cl$	$100D$	$1D$	$50D$	$100D$	$2D$	$50D$
Baseline	- Random	298.24 \pm 0.23	4.66 \pm 0.03	200.8 \pm 0.6	73.01 \pm 0.17	1704.3 \pm 9.3	50.0 \pm 0.1	50.0 \pm 0.3	49.9 \pm 0.3
	- Optimization	100.88 \pm 0.00	0.43 \pm 0.00	20.1 \pm 0.0	0.36 \pm 0.00	309.2 \pm 0.2	71.2 \pm 0.0	96.8 \pm 0.0	76.1 \pm 0.0
	- Langevin	101.92 \pm 0.04	0.44 \pm 0.00	21.8 \pm 1.0	0.31 \pm 0.00	N/A	65.5 \pm 0.5	96.1 \pm 0.0	70.1 \pm 0.2
	- HMC	102.01 \pm 0.01	0.46 \pm 0.01	17.8 \pm 0.1	0.38 \pm 0.01	303.9 \pm 2.5	62.6 \pm 0.2	91.7 \pm 0.2	68.0 \pm 0.4
Fwd-KL	GRU	133.22 \pm 0.58	2.36 \pm 0.02	139.4 \pm 1.0	51.45 \pm 0.03	1346.5 \pm 6.8	57.9 \pm 0.2	59.6 \pm 0.2	58.6 \pm 0.2
	DeepSets	129.69 \pm 0.74	2.35 \pm 0.02	149.8 \pm 0.8	51.90 \pm 1.54	1357.5 \pm 5.3	50.8 \pm 0.1	49.9 \pm 0.3	49.9 \pm 0.3
	Transformer	108.98 \pm 0.10	2.40 \pm 0.02	64.3 \pm 3.7	50.81 \pm 0.53	1319.9 \pm 12.2	62.4 \pm 0.0	59.9 \pm 0.2	58.8 \pm 0.2
Rev-KL	GRU	105.14 \pm 0.10	0.46 \pm 0.01	62.6 \pm 0.1	2.31 \pm 0.13	316.3 \pm 6.2	63.4 \pm 0.2	88.8 \pm 0.5	68.4 \pm 0.3
	DeepSets	105.06 \pm 0.21	0.48 \pm 0.02	64.1 \pm 0.2	0.98 \pm 0.12	451.9 \pm 2.8	61.3 \pm 0.1	91.0 \pm 0.5	61.7 \pm 0.1
	Transformer	104.71 \pm 0.12	0.47 \pm 0.01	32.0 \pm 0.5	0.81 \pm 0.02	278.3 \pm 1.1	67.7 \pm 0.1	90.0 \pm 0.2	73.7 \pm 0.3
Fwd-KL	GRU	125.84 \pm 1.98	0.60 \pm 0.07	138.3 \pm 1.0	38.41 \pm 0.36	1097.4 \pm 9.5	58.0 \pm 0.1	61.2 \pm 0.8	60.2 \pm 0.1
	DeepSets	133.23 \pm 1.93	0.58 \pm 0.03	153.2 \pm 0.8	43.31 \pm 2.06	1120.0 \pm 5.5	50.5 \pm 0.1	49.6 \pm 0.2	50.1 \pm 0.1
	Transformer	108.48 \pm 0.16	0.59 \pm 0.08	63.1 \pm 2.0	39.70 \pm 0.52	1073.3 \pm 1.5	63.6 \pm 0.1	60.9 \pm 0.3	60.3 \pm 0.1
Rev-KL	GRU	105.19 \pm 0.03	0.47 \pm 0.01	71.3 \pm 1.3	2.31 \pm 0.41	302.9 \pm 5.6	63.4 \pm 0.1	90.4 \pm 0.7	66.2 \pm 0.1
	DeepSets	105.09 \pm 0.06	0.49 \pm 0.01	76.8 \pm 1.8	0.83 \pm 0.02	454.1 \pm 10.2	59.1 \pm 0.5	89.1 \pm 0.3	62.9 \pm 0.1
	Transformer	104.91 \pm 0.11	0.46 \pm 0.00	33.1 \pm 0.3	0.99 \pm 0.07	274.0 \pm 1.3	68.1 \pm 0.2	91.1 \pm 0.2	72.6 \pm 0.1

Table 2: **Variable-Dimension Posterior Prediction.** Results for estimating the mean of a Gaussian (Gaussian), means of a Gaussian mixture model (GMM), (non-)linear regression (NLR/LR) and (non-)linear binary classification (NLC/LC). In this setting, for each task only a single model was trained to solve for different dimensional problems; e.g. the same model is used to evaluate $1D$ and $50D$ nonlinear regression.

We note that the form of the solution remains the same irrespective of whether the underlying problem is a 2- or 100-dimensional Bayesian linear regression problem. Additionally, we can see that a low-dimensional problem can just be embedded into a high-dimensional space, with the extra features and parameters set to 0, akin to the procedure of masking unnecessary dimensions. This simple but strong insight allows us to amortize q_φ over datasets with varying dimensionalities by embedding them in a higher dimensional space.

3.2 Handling Model Misspecification

Knowledge about the true likelihood model underlying the data-generating scheme is often unknown. Practitioners thus rely on making a modeling assumption by prescribing a particular likelihood model for the data and then fitting its parameters to best explain the data. For example, the underlying model of whether an email is spam or not may be unknown, but one can still assume a linear model to solve this classification problem. This defines the basis of model misspecification, where there is a mismatch between the underlying true model and the assumption used.

In contrast to SBI-NPE approaches, our method can efficiently handle model misspecification. This can be seen from the contrast between Equations 9 and 12, where the former can only be trained when χ is defined according to the modeling choice while the latter doesn't put any constraints on χ . Thus, in the presence of a stream of real-world data, SBI-NPE cannot leverage this data during training because the corresponding parameters in the assumed likelihood model are not known, while our method can leverage it during training and can thus lead to more robust predictions on data that is of actual practical significance.

A simple example of this setup is when the true data might be coming from a Gaussian Process, but we might not know this. To still be able to efficiently make predictions, we can model it with a linear relation. However, SBI-NPE approaches can only train the amortized posterior for the linear model using linear data,

$\chi_{real} (\rightarrow)$	q_φ	<i>Data</i>	Linear	MLP Nonlinear	GP Nonlinear
$\chi_{sim} (\rightarrow)$		<i>Model</i>	NLR	LR	NLR
Fwd-KL	Gaussian		15.454 \pm 0.246	2.216 \pm 0.097	14.733 \pm 0.513
Rev-KL			0.382 \pm 0.003	1.892 \pm 0.113	0.155 \pm 0.006
+ <i>switched data</i>			0.367 \pm 0.006	1.226 \pm 0.001	0.066 \pm 0.004
Fwd-KL	Flow		7.949 \pm 0.419	1.632 \pm 0.070	8.557 \pm 0.561
Rev-KL			0.347 \pm 0.001	1.471 \pm 0.016	0.120 \pm 0.005
+ <i>switched data</i>			0.346 \pm 0.002	1.226 \pm 0.004	0.055 \pm 0.002

Table 3: **Model Misspecification:** We see benefits of our proposed approach when the true underlying data generating process is not known. χ_{real} denotes the real data-generating distribution, χ_{sim} the simulated one according to the (wrongly) assumed model. We train different models on χ_{sim} , with switched data denoting training on χ_{real} . Evaluation is done solely on χ_{real} .

while our method can actually leverage the GP data that we observe. This leads to better generalization to linear modeling of GP data and can be further leveraged in other model misspecification setups.

4 Experiments

To showcase the wide applicability of our approach, we perform experiments on various well-known probabilistic models encompassing supervised and unsupervised scenarios. In particular, we look at the problem of estimating the Bayesian posterior over the (a) mean of a Gaussian distribution (GM), (b) means of a Gaussian mixture model (GMM), (c) parameters of a (non-)linear regression model (NLR/LR), and (d) parameters of a (non-)linear classification model (NLC/LC). We refer the readers to Appendix C for particulars about the probabilistic models, including their likelihoods and priors considered. Throughout our experiments, we observe superior performance of our proposed approach in terms of predictive performance metrics, especially in problems with high-dimensional and multi-modal posteriors.

We generally consider two types of baselines: dataset-specific and amortized. For dataset-specific baselines, we use the prior (Random), perform maximum likelihood estimation using gradient-based optimization (Optimization) as well as an approximate Bayesian inference procedure through Langevin and Hamiltonian based MCMC sampling, which also uses the gradient information. Such baselines rely on iterative procedures and must be run independently for different datasets. For amortized baselines, we consider the SBI-NPE based forward-KL baseline as outlined in Equation 9. We refer the readers to Appendices F, D, and H for details about the experiments, metrics, and additional results, respectively.

Zero-Shot Posterior Approximation.

Given a known probabilistic model $p(\mathcal{D}, \theta)$, we train an amortized inference system $q_\varphi(\cdot|\mathcal{D})$ to approximate the often intractable posterior over the parameter θ for different datasets with varying cardinality. Figure 1 visualizes the performance of the amortized model on different probabilistic models spanning supervised and unsupervised learning for low dimensional problems. We see that the proposed approach provides reliable samples of the parameters for previously unseen datasets, showing that it can indeed generalize.

Next, we empirically evaluate the performance of our proposed approach in more complex, high-dimensional setups for the same set of probabilistic models. We set the data-generating distribution χ to be according to the probabilistic model in order to allow comparisons with SBI-NPE, which can only be trained in this particular setting. Table 1 highlights the performance of our proposed approach in high-dimensional frameworks, showcasing its superior performance to SBI-NPE benchmarks, as well as comparative to dataset-specific non-amortized baselines. For multiple high-dimensional problems, we see that the proposed amortized posterior performs similarly to Optimization and MCMC without being slow and iterative at inference.

Generalizing to Variable Feature Dimensions

Objective	q_φ	Model	L_2 Loss (\downarrow)		Accuracy (\uparrow)	
			LR	NLR	LC	NLC
Baseline	-	Random	23.52 \pm 0.42	209.35 \pm 9.92	50.10 \pm 0.17	50.84 \pm 1.00
Fwd-KL	Gaussian	GRU	8.95 \pm 0.47	84.63 \pm 3.96	76.34 \pm 1.58	60.00 \pm 0.98
		DeepSets	10.81 \pm 0.08	97.51 \pm 2.96	68.26 \pm 0.31	50.84 \pm 0.99
		Transformer	9.35 \pm 0.99	111.07 \pm 6.55	63.97 \pm 3.21	60.39 \pm 0.53
Rev-KL	Gaussian	GRU	9.78 \pm 2.75	17.01 \pm 5.16	79.71 \pm 1.12	77.19 \pm 0.21
		DeepSets	9.26 \pm 0.10	8.03 \pm 0.23	77.40 \pm 0.14	72.63 \pm 0.14
		Transformer	7.30 \pm 0.21	8.23 \pm 1.26	76.96 \pm 1.58	71.33 \pm 5.44
Fwd-KL	Flow	GRU	8.28 \pm 0.33	52.04 \pm 3.39	76.45 \pm 1.43	61.10 \pm 0.61
		DeepSets	12.94 \pm 0.41	71.69 \pm 2.38	67.99 \pm 2.15	49.74 \pm 0.76
		Transformer	9.64 \pm 0.50	84.45 \pm 7.88	68.26 \pm 3.23	61.65 \pm 1.22
Rev-KL	Flow	GRU	8.17 \pm 0.25	17.35 \pm 8.02	66.30 \pm 2.81	78.55 \pm 1.27
		DeepSets	11.05 \pm 0.35	8.54 \pm 0.49	78.18 \pm 0.13	71.69 \pm 0.18
		Transformer	7.48 \pm 0.26	8.80 \pm 1.80	72.75 \pm 5.73	78.23 \pm 0.89

Table 4: **Tabular Experiments:** Zero-shot performance of the amortized variable-dimensional models across (non-)linear regression and classification real-world tabular tasks. All the models are trained solely on simulated data, and evaluated zero-shot on real-world data with varying number of both features and training (observations that are fed as in-context amortization) observations.

Next, we turn our attention to performing amortization additionally over different datasets *with a variable number of feature dimensions*, as highlighted in Section 3.1. To do this, we embed all low-dimensional problems in a 100-dimensional space by masking the unnecessary dimensions, and then train an amortized VI system. A key benefit of this approach is that instead of having to train different models for different dimensional problems, we can rely on a single model that can generalize across all. We visualize samples from the approximate posterior for low-dimensional problems in Figure 2 and further provide quantitative predictive metrics for high-dimensional problems in Table 2. Our results highlight that the proposed approach can generalize quite well to variable dimensional setups without requiring an iterative procedure at inference.

Handling Model Misspecification

As discussed in Section 3 and prior work (Müller et al., 2021; Hollmann et al., 2022), the performance and use-case of such amortized systems relies heavily on the choice of the data-generating distribution χ . So far, we only look at cases where χ defines simulations according to the probabilistic model to allow for a fair comparison with SBI-NPE. This setup, however, implicitly assumes that we know the ground-truth data-generating process, which is often not the case. A more common scenario is where we do not know the data-generating process but have access to samples from it defined through χ_{real} . To make predictions for new query points, we, as statisticians, assume a probabilistic model $p(\mathbf{x}|\theta)$ to explain the data, which may be quite far from the true underlying model.

We first note that χ_{real} cannot be used to train a forward KL model (SBI-NPE) and is thus unusable in this case. To resolve this issue, let χ_{sim} denote the data generating distribution obtained from simulations via the assumed model, which can then be used to train the forward KL model. We compare this setup with our proposed approach, which can be trained on either χ_{sim} or χ_{real} . However, for this setup, we always perform evaluation on the real data of interest, i.e. χ_{real} .

For our experiments, we pick χ_{real} and χ_{sim} by simulation through probabilistic models corresponding to linear, MLP, and Gaussian Process regression, such that χ_{real} and χ_{sim} are never the same. Table 3 showcases experimental results of model misspecification, where we see that the forward KL models suffer a lot due to mismatch between the training and evaluation data, compounded by the potential overestimation of variance (mode-averaging) while our proposed method significantly outperforms it. Even further, we see that the proposed method can be directly trained on χ_{real} (switched data), leading to even better performance – for eg., when χ_{real} denotes data generated from a Gaussian Process and the probabilistic model assumed

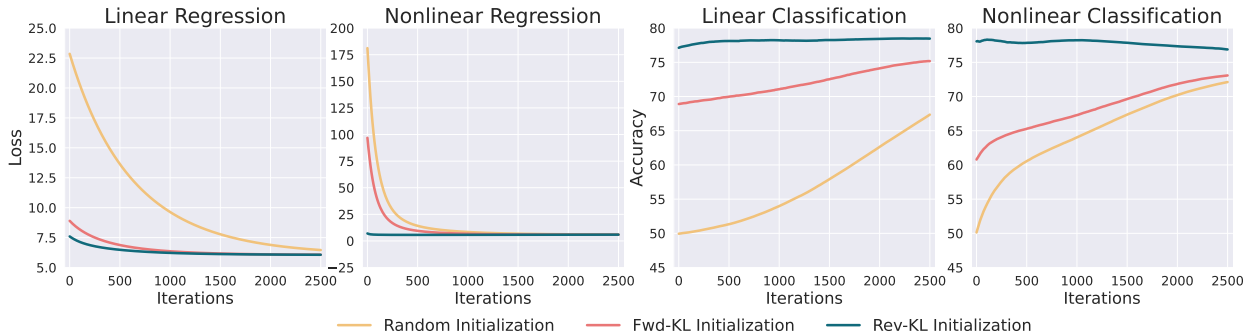


Figure 3: **Tabular Experiments:** Initializing parameters from the proposed amortized model leads to good zero-shot performance and often optimal initialization across (non-)linear regression and classification tasks. In the MLP-based (last column), our proposed approach outperforms forward KL when trained on the same simulated data from the MLP, and does even better when trained on the GP-based data. We refer the readers to Appendix F.3 and H.3 for experimental details and additional results on model misspecification.

Applications to Tabular Domains

While we see clear benefits in simulated settings, we next turn our attention to real-world scenarios and consider a suite of tasks from the OpenML platform for both regression and binary classification by filtering them out from the *OpenML-CTR23 - A curated tabular regression benchmarking suite* (Fischer et al., 2023) for regression and *OpenML-CC18 Curated Classification benchmark* (Bischl et al., 2019) for classification problems (details in Appendix G). We end up with 9 regression and 13 classification datasets with varying number of features and use our amortized inference systems trained on simulated data with variable feature dimensions to predict the parameters of interest.

After initializing from the inference model, we further train the parameters of the respective probabilistic models with maximum a-posteriori (*map*) objective and compare the performance with a corresponding model initialized from the prior. Even in this extreme case of domain shift, we see in Table 4 that the amortized model provides good initializations for different real-world tasks zero-shot after training only on simulated data. Figure 3 also shows that the amortized model surpasses the SBI-NPE baseline in terms of performance and convergence speed for both linear and nonlinear setups. While Figure 3 only provides a normalized aggregated performance over all datasets considered for a q_φ parameterized as a diagonal Gaussian, we refer the readers to Appendix H.4 for results on individual datasets with different q_φ , as well as Appendix F.4 for implementation details.

Evaluating Posterior Inference

While comparing with the true posterior is hard due to its intractability, it is still available for the problem of estimating the mean of a Gaussian distribution as well as for Bayesian Linear Regression. We plot the kernel density estimate of the samples obtained from the true posterior, the amortized forward, and the reverse KL model in Figure 4 (Right), which shows that both the amortization setups efficiently capture the true posterior distribution. Further, we quantify the quality of the approximate posteriors through a symmetric KL Divergence metric in Table 5, showing that the amortized model obtains good approximations but does worsen with increasing dimensionality, with reverse KL modeling outperforming SBI-NPE on high dimensional linear regression.

q_φ	Model	Symmetric KL Divergence (\downarrow)			
		Gaussian Mean		LR	
		2D	100D	1D	100D
Baseline -	Random	44.32	46.78	179.2	186.8
Fwd-KL	GRU	0.018 \pm 0.007	0.07 \pm 0.00	0.04 \pm 0.01	81.8 \pm 0.2
	DeepSets	0.037 \pm 0.015	0.22 \pm 0.01	0.06 \pm 0.00	82.0 \pm 0.4
	Transformer	0.030 \pm 0.008	0.06 \pm 0.00	0.04 \pm 0.01	20.6 \pm 1.0
Rev-KL	GRU	0.017 \pm 0.005	0.08 \pm 0.00	0.03 \pm 0.00	78.7 \pm 1.5
	DeepSets	0.029 \pm 0.002	0.19 \pm 0.01	0.05 \pm 0.00	104.5 \pm 8.8
	Transformer	0.035 \pm 0.013	0.05 \pm 0.00	0.03 \pm 0.00	32.7 \pm 1.0

Table 5: **Normalized Symmetric KL Divergence.** Amortized models approximate the true posterior well in tasks with tractable posteriors, when compared to the prior .

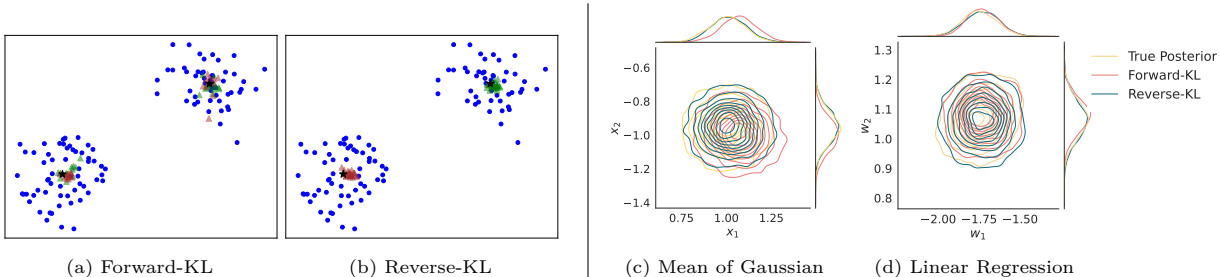


Figure 4: **Left**: Estimation of the means of a GMM, where red and green samples denote the first and second mean vectors. Unlike in reverse KL, the cluster labels switch in forward KL, highlighting its ability to capture underlying multi-modality. **Right**: Kernel density estimation of the true posterior, overlaid with estimates from forward and reverse KL systems, for different probabilistic models.

For more complex problems, we additionally compute the squared wasserstein metric \mathcal{W}_2^2 between samples from the amortized approximate posterior and from multiple chains of Langevin MCMC. Table 6 shows that reverse KL approaches do slightly better in high-dimensional setups, while for low-dimensional multi-modal scenarios (eg. GMM), forward KL approaches fare better; highlighting that they suffer less from mode collapse. However, it is imperative to know that this metric only provides a crude proxy to the quality of the posterior, since MCMC methods only provide guarantees of true samples asymptotically.

5 Discussion

Having presented a method to amortize posterior estimation for different probabilistic models and the corresponding empirical evidence, we now highlight some of our key findings.

Posterior Inference. Our results indicate that the proposed amortized Bayesian posterior estimate generalizes quite well zero-shot to novel tasks without the need for fine-tuning. While we see that in low-dimensional multi-modal settings, forward KL performs better at capturing multiple modes, it fails to capture complex posteriors in high dimensions. Our posterior predictive results additionally indicate that forward KL often performs comparable to random chance (prior) in high-dimensional, multi-modal setups. In contrast, reverse KL provides reasonable estimates, alluding to its ability to better capture at least a mode of the posterior and highlights its superiority and usability when modeling predictions.

Architecture Choices. To better understand the architectural inductive biases that lead to better posterior estimates, we compare permutation invariant architectures like DeepSets and Transformers with non invariant architecture like GRU. We first note that GRU outperforms DeepSets even though the latter is permutation invariant, which can be explained by the fact that the former can still learn invariant structures after seeing a lot of datasets during training while the latter is restrictive in its invariance due to the choice of aggregation operator, which we keep as mean in our setting. On the other hand, we see that Transformer outperforms both DeepSets and GRUs as it does not rely on a fixed aggregation scheme (e.g. sum or mean pooling) while still respecting the invariant structure of the posterior and learning it in a context-aware fashion. However, we do see that in some rare cases of OoD generalization, DeepSets and GRUs can still outperform Transformers (Appendix H.3).

Capacity of q_φ . Next, we see that increasing the capacity of q_φ with normalizing flows only helps marginally for the reverse KL objective but substantially for the forward KL setup. We hypothesize that given the mode-seeking tendency of reverse KL, even with the capacity to model different modes, the algorithm seeks and latches to only a single mode and capturing multiple modes in this setup is challenging, whereas in forward KL setup without additional capacity the model overestimates the variance a lot.

Forward vs Reverse KL. Our experiments on GMM show that forward KL objective leads to the learning of a multi-modal distribution while reverse KL only captures one mode (Figure 4, **Left**). However, in high-dimensional multi-modal settings like BNNs, the former does not lead to learning of a reasonable distribution as it attempts to cover all the modes while the latter does not cover multiple modes but better models an individual mode (Tables 1 and 2; Appendix H.1 and H.2). Furthermore, unlike forward KL, the reverse KL

Objective q_φ	Model	$\mathcal{W}_2^2 (\downarrow)$								
		Gaussian	GMM	LR	NLR		LC	NLC		
		100D	5D 2 cl	100D	1D	25D	100D	2D	25D	
Baseline	-	Random	13.96±0.00	4.01±0.00	13.53±0.00	11.21±0.00	36.46±0.00	16.72±0.00	14.71±0.00	36.67±0.00
Fwd-KL	Gaussian	GRU	1.37±0.00	2.35±0.01	10.42±0.02	11.10±0.00	36.33±0.01	15.25±0.01	14.67±0.00	36.66±0.00
		DeepSets	1.55±0.01	2.35±0.01	10.39±0.03	11.11±0.01	36.41±0.01	16.72±0.00	14.68±0.00	36.66±0.00
		Transformer	1.41±0.00	2.40±0.01	5.83±0.09	11.09±0.01	36.32±0.01	14.70±0.01	14.67±0.00	36.66±0.00
Rev-KL	Gaussian	GRU	1.34±0.00	2.98±0.01	7.39±0.03	11.31±0.01	35.92±0.32	12.31±0.01	14.14±0.01	35.16±0.00
		DeepSets	1.38±0.02	2.98±0.02	7.58±0.05	11.31±0.02	35.58±0.20	12.93±0.06	14.10±0.01	35.05±0.00
		Transformer	1.34±0.01	2.98±0.03	4.84±0.03	11.38±0.01	35.86±0.08	12.84±0.04	14.17±0.02	35.36±0.00
Fwd-KL	Flow	GRU	1.37±0.00	1.71±0.16	10.18±0.04	11.09±0.02	36.35±0.01	15.29±0.02	14.66±0.01	36.66±0.00
		DeepSets	1.58±0.01	1.81±0.08	10.48±0.14	11.08±0.01	36.41±0.00	16.72±0.00	14.67±0.01	36.66±0.00
		Transformer	1.40±0.01	1.20±0.26	5.64±0.23	11.08±0.00	36.32±0.02	14.66±0.01	14.65±0.01	36.66±0.00
Rev-KL	Flow	GRU	1.33±0.00	2.99±0.02	7.48±0.06	11.15±0.04	35.97±0.04	13.51±0.02	14.33±0.01	35.79±0.01
		DeepSets	1.41±0.02	2.96±0.01	8.40±0.09	11.15±0.04	36.02±0.09	13.61±0.05	14.32±0.01	35.69±0.01
		Transformer	1.33±0.01	3.00±0.04	4.90±0.15	11.21±0.02	36.00±0.17	13.63±0.02	14.37±0.01	35.88±0.02

Table 6: **Sample Based Metrics.** We compute the squared 2–Wasserstein metric between samples from the approximate posterior and samples obtained from Langevin MCMC algorithm, over a range of different tasks.

paradigm can be trained without observing θ but does require a known differentiable likelihood. We show an application of this setting where we don’t have access to the ground-truth model but only to samples from it.

6 Conclusion

We show that Bayesian posterior inference can be amortized for a broad class of probabilistic models and explore a variety of design decisions associated with it. In particular, we show that reverse KL is effective for learning the amortization network and has significant benefits when modeling the predictive, especially in the presence of model misspecification and generalization to out-of-domain real-world setups. It provides an exciting direction of research which could reduce the load of real-world, complex, and iterative approximations through quick and cheap inference over a trained amortized network. Even further, we believe that our approach provides a direction into learning of generalist in-context Bayesian estimator capable of estimating the posterior for novel datasets zero-shot.

We believe that scaling our approach to more complex probabilistic models as well as leveraging existing real-world data in the training paradigm can lead to even better amortized posterior estimators. Another important future work is to combine our approach with diffusion-based (Zhang & Chen, 2021; Vargas et al., 2023) and other off-policy variational systems (Bengio et al., 2021) for modeling more complex high dimensional distributions.

Impact Statement

This work studies amortizing variational inference for Bayesian posterior estimation which is a widespread strategy for performing inference in statistics. It provides a natural way of quantifying uncertainty and potentially leading to more robust predictions. While we do not foresee any negative impacts of progress in this area, we encourage caution when applying the methodologies in practice.

References

- Akhound-Sadegh, T., Rector-Brooks, J., Bose, A. J., Mittal, S., Lemos, P., Liu, C.-H., Sendera, M., Ravanbakhsh, S., Gidel, G., Bengio, Y., et al. Iterated denoising energy matching for sampling from boltzmann densities. *arXiv preprint arXiv:2402.06121*, 2024.
- Ardizzone, L., Bungert, T., Draxler, F., Köthe, U., Kruse, J., Schmier, R., and Sorrenson, P. FrEIA: Framework for easily invertible architectures, 2018. URL <https://github.com/vislearn/FrEIA>.
- Ardizzone, L., Bungert, T., Draxler, F., Köthe, U., Kruse, J., Schmier, R., and Sorrenson, P. Framework for Easily Invertible Architectures (FrEIA), 2018-2022. URL <https://github.com/vislearn/FrEIA>.
- Arenz, O., Dahlinger, P., Ye, Z., Volpp, M., and Neumann, G. A unified perspective on natural gradient variational inference with gaussian mixture models. *arXiv preprint arXiv:2209.11533*, 2022.
- Athanasopoulos, G., Hyndman, R. J., Kourentzes, N., and Panagiotelis, A. Forecast reconciliation: A review. *International Journal of Forecasting*, 2023. ISSN 0169-2070. doi: <https://doi.org/10.1016/j.ijforecast.2023.10.010>. URL <https://www.sciencedirect.com/science/article/pii/S0169207023001097>.
- Bengio, E., Jain, M., Korablyov, M., Precup, D., and Bengio, Y. Flow network based generative models for non-iterative diverse candidate generation. *Advances in Neural Information Processing Systems*, 34: 27381–27394, 2021.
- Berner, J., Richter, L., and Ullrich, K. An optimal control perspective on diffusion-based generative modeling. *arXiv preprint arXiv:2211.01364*, 2022.
- Bingham, E., Chen, J. P., Jankowiak, M., Obermeyer, F., Pradhan, N., Karaletsos, T., Singh, R., Szerlip, P., Horsfall, P., and Goodman, N. D. Pyro: Deep universal probabilistic programming. *The Journal of Machine Learning Research*, 20(1):973–978, 2019.
- Bischl, B., Casalicchio, G., Feurer, M., Hutter, F., Lang, M., Mantovani, R. G., van Rijn, J. N., and Vanschoren, J. Openml benchmarking suites. *arXiv:1708.03731v2 [stat.ML]*, 2019.
- Bishop, C. M. and Nasrabadi, N. M. *Pattern recognition and machine learning*, volume 4. Springer, 2006.
- Bitzer, M., Meister, M., and Zimmer, C. Amortized inference for gaussian process hyperparameters of structured kernels. *arXiv preprint arXiv:2306.09819*, 2023.
- Blei, D. M., Kucukelbir, A., and McAuliffe, J. D. Variational inference: A review for statisticians. *Journal of the American statistical Association*, 112(518):859–877, 2017.
- Brooks, S. Markov chain monte carlo method and its application. *Journal of the royal statistical society: series D (the Statistician)*, 47(1):69–100, 1998.
- Carpenter, B., Gelman, A., Hoffman, M. D., Lee, D., Goodrich, B., Betancourt, M., Brubaker, M. A., Guo, J., Li, P., and Riddell, A. Stan: A probabilistic programming language. *Journal of statistical software*, 76, 2017.
- Chauhan, V. K., Zhou, J., Lu, P., Molaei, S., and Clifton, D. A. A brief review of hypernetworks in deep learning. *arXiv preprint arxiv:2306.06955*, 2023.
- Chen, Y., Garnett, R., and Montgomery, J. M. Polls, context, and time: A dynamic hierarchical bayesian forecasting model for us senate elections. *Political Analysis*, 31(1):113–133, 2023. doi: 10.1017/pan.2021.42.
- Cooper, I., Mondal, A., and Antonopoulos, C. G. A sir model assumption for the spread of covid-19 in different communities. *Chaos, Solitons & Fractals*, 139:110057, 2020. ISSN 0960-0779. doi: <https://doi.org/10.1016/j.chaos.2020.110057>. URL <https://www.sciencedirect.com/science/article/pii/S0960077920304549>.

- Cranmer, K., Brehmer, J., and Louppe, G. The frontier of simulation-based inference. *Proceedings of the National Academy of Sciences*, 117(48):30055–30062, May 2020. ISSN 1091-6490. doi: 10.1073/pnas.1912789117. URL <http://dx.doi.org/10.1073/pnas.1912789117>.
- Dinh, L., Sohl-Dickstein, J., and Bengio, S. Density estimation using real NVP. *5th International Conference on Learning Representations, ICLR 2017 - Conference Track Proceedings*, 2017. URL <http://arxiv.org/abs/1605.08803>.
- Finn, C., Abbeel, P., and Levine, S. Model-agnostic meta-learning for fast adaptation of deep networks. In *International conference on machine learning*, pp. 1126–1135. PMLR, 2017.
- Fischer, S. F., Feurer, M., and Bischl, B. OpenML-CTR23 – a curated tabular regression benchmarking suite. In *AutoML Conference 2023 (Workshop)*, 2023. URL <https://openreview.net/forum?id=HebA0oMm94>.
- Gardner, J., Pleiss, G., Weinberger, K. Q., Bindel, D., and Wilson, A. G. Gpytorch: Blackbox matrix-matrix gaussian process inference with gpu acceleration. *Advances in neural information processing systems*, 31, 2018.
- Garnelo, M., Rosenbaum, D., Maddison, C., Ramalho, T., Saxton, D., Shanahan, M., Teh, Y. W., Rezende, D., and Eslami, S. A. Conditional neural processes. In *International conference on machine learning*, pp. 1704–1713. PMLR, 2018a.
- Garnelo, M., Schwarz, J., Rosenbaum, D., Viola, F., Rezende, D. J., Eslami, S., and Teh, Y. W. Neural processes. *arXiv preprint arXiv:1807.01622*, 2018b.
- Geffner, T., Papamakarios, G., and Mnih, A. Compositional score modeling for simulation-based inference. 2023.
- Gelfand, A. E. Gibbs sampling. *Journal of the American statistical Association*, 95(452):1300–1304, 2000.
- Gelman, A., Carlin, J. B., Stern, H. S., Dunson, D. B., Vehtari, A., and Rubin, D. B. *Bayesian Data Analysis, Third Edition*. CRC Press, November 2013. ISBN 9781439840955. URL <https://play.google.com/store/books/details?id=ZXL6AQAQBAJ>.
- Glaeser, R. M., Nogales, E., and Chiu, W. *Single-particle Cryo-EM of Biological Macromolecules*. 2053-2563. IOP Publishing, 2021. ISBN 978-0-7503-3039-8. doi: 10.1088/978-0-7503-3039-8. URL <https://dx.doi.org/10.1088/978-0-7503-3039-8>.
- Gordon, J., Bruinsma, W. P., Foong, A. Y., Requeima, J., Dubois, Y., and Turner, R. E. Convolutional conditional neural processes. *arXiv preprint arXiv:1910.13556*, 2019.
- Grant, E., Finn, C., Levine, S., Darrell, T., and Griffiths, T. Recasting gradient-based meta-learning as hierarchical bayes. *arXiv preprint arXiv:1801.08930*, 2018.
- Higgins, I., Matthey, L., Pal, A., Burgess, C., Glorot, X., Botvinick, M., Mohamed, S., and Lerchner, A. beta-VAE: Learning basic visual concepts with a constrained variational framework. In *International Conference on Learning Representations*, 2017. URL <https://openreview.net/forum?id=Sy2fzU9gl>.
- Hoffman, M. D., Blei, D. M., Wang, C., and Paisley, J. Stochastic variational inference. *Journal of Machine Learning Research*, 2013.
- Hoffman, M. D., Gelman, A., et al. The no-u-turn sampler: adaptively setting path lengths in hamiltonian monte carlo. *J. Mach. Learn. Res.*, 15(1):1593–1623, 2014.
- Hollmann, N., Müller, S., Eggensperger, K., and Hutter, F. TabPFN: A transformer that solves small tabular classification problems in a second. *arXiv preprint arXiv:2207.01848*, 2022.
- Hospedales, T., Antoniou, A., Micaelli, P., and Storkey, A. Meta-learning in neural networks: A survey. *IEEE Transactions on Pattern Analysis & Machine Intelligence*, 44(09):5149–5169, sep 2022. ISSN 1939-3539. doi: 10.1109/TPAMI.2021.3079209.

- Kim, H., Mnih, A., Schwarz, J., Garnelo, M., Eslami, A., Rosenbaum, D., Vinyals, O., and Teh, Y. W. Attentive neural processes. *arXiv preprint arXiv:1901.05761*, 2019.
- Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Kingma, D. P. and Dhariwal, P. Glow: Generative flow with invertible 1x1 convolutions. *Advances in neural information processing systems*, 31, 2018.
- Kingma, D. P. and Welling, M. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- Kingma, D. P., Welling, M., et al. An introduction to variational autoencoders. *Foundations and Trends® in Machine Learning*, 12(4):307–392, 2019.
- Kobyzev, I., Prince, S. J., and Brubaker, M. A. Normalizing flows: An introduction and review of current methods. *IEEE transactions on pattern analysis and machine intelligence*, 43(11):3964–3979, 2020.
- Koch, G., Zemel, R., Salakhutdinov, R., et al. Siamese neural networks for one-shot image recognition. In *ICML deep learning workshop*, volume 2. Lille, 2015.
- Krueger, D., Huang, C.-W., Islam, R., Turner, R., Lacoste, A., and Courville, A. Bayesian hypernetworks. *arXiv preprint arxiv:1710.04759*, 2017.
- Lee, J., Lee, Y., Kim, J., Kosiorek, A., Choi, S., and Teh, Y. W. Set transformer: A framework for attention-based permutation-invariant neural networks. In *International conference on machine learning*, pp. 3744–3753. PMLR, 2019.
- Lin, W., Schmidt, M., and Khan, M. E. Handling the positive-definite constraint in the bayesian learning rule. In *International conference on machine learning*, pp. 6116–6126. PMLR, 2020.
- Liu, S., Sun, X., Ramadge, P. J., and Adams, R. P. Task-agnostic amortized inference of gaussian process hyperparameters. *Advances in Neural Information Processing Systems*, 33:21440–21452, 2020.
- Lorch, L., Sussex, S., Rothfuss, J., Krause, A., and Schölkopf, B. Amortized inference for causal structure learning. *Advances in Neural Information Processing Systems*, 35:13104–13118, 2022.
- Morris, Q. Recognition networks for approximate inference in bn20 networks. *arXiv preprint arXiv:1301.2295*, 2013.
- Müller, S., Hollmann, N., Arango, S. P., Grabocka, J., and Hutter, F. Transformers can do bayesian inference. *arXiv preprint arXiv:2112.10510*, 2021.
- Paige, B. and Wood, F. Inference networks for sequential monte carlo in graphical models. In *International Conference on Machine Learning*, pp. 3040–3049. PMLR, 2016.
- Pakman, A., Wang, Y., Mitelut, C., Lee, J., and Paninski, L. Neural clustering processes. In *International Conference on Machine Learning*, pp. 7455–7465. PMLR, 2020.
- Papamakarios, G., Nalisnick, E., Rezende, D. J., Mohamed, S., and Lakshminarayanan, B. Normalizing flows for probabilistic modeling and inference. *The Journal of Machine Learning Research*, 22(1):2617–2680, 2021.
- Radev, S. T., Mertens, U. K., Voss, A., Ardizzone, L., and Köthe, U. Bayesflow: Learning complex stochastic models with invertible neural networks. *IEEE transactions on neural networks and learning systems*, 33(4): 1452–1466, 2020.
- Rezende, D. and Mohamed, S. Variational inference with normalizing flows. In *International conference on machine learning*, pp. 1530–1538. PMLR, 2015.
- Rezende, D. J., Mohamed, S., and Wierstra, D. Stochastic backpropagation and approximate inference in deep generative models. In *International conference on machine learning*, pp. 1278–1286. PMLR, 2014.

- Simpson, F., Davies, I., Lalchand, V., Vullo, A., Durrande, N., and Rasmussen, C. E. Kernel identification through transformers. *Advances in Neural Information Processing Systems*, 34:10483–10495, 2021.
- Song, Y., Sohl-Dickstein, J., Kingma, D. P., Kumar, A., Ermon, S., and Poole, B. Score-based generative modeling through stochastic differential equations. *arXiv preprint arXiv:2011.13456*, 2020.
- Stuhlmüller, A., Taylor, J., and Goodman, N. Learning stochastic inverses. *Advances in neural information processing systems*, 26, 2013.
- Sun, Z., Ozay, M., and Okatani, T. Hypernetworks with statistical filtering for defending adversarial examples. *arXiv preprint arxiv:1711.01791*, 2017.
- Sung, F., Yang, Y., Zhang, L., Xiang, T., Torr, P. H., and Hospedales, T. M. Learning to compare: Relation network for few-shot learning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1199–1208, 2018.
- Tang, L., Zhou, Y., Wang, L., Purkayastha, S., Zhang, L., He, J., Wang, F., and Song, P. X.-K. A review of multi-compartment infectious disease models. *International Statistical Review*, 88(2):462–513, 2020. doi: <https://doi.org/10.1111/insr.12402>. URL <https://onlinelibrary.wiley.com/doi/abs/10.1111/insr.12402>.
- Vargas, F., Grathwohl, W., and Doucet, A. Denoising diffusion samplers. *arXiv preprint arXiv:2302.13834*, 2023.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., and Polosukhin, I. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- Vinyals, O., Blundell, C., Lillicrap, T., Wierstra, D., et al. Matching networks for one shot learning. *Advances in neural information processing systems*, 29, 2016.
- Von Oswald, J., Niklasson, E., Randazzo, E., Sacramento, J., Mordvintsev, A., Zhmoginov, A., and Vladymyrov, M. Transformers learn in-context by gradient descent. In *International Conference on Machine Learning*, pp. 35151–35174. PMLR, 2023.
- von Oswald, J., Niklasson, E., Schlegel, M., Kobayashi, S., Zucchet, N., Scherrer, N., Miller, N., Sandler, M., Vladymyrov, M., Pascanu, R., et al. Uncovering mesa-optimization algorithms in transformers. *arXiv preprint arXiv:2309.05858*, 2023.
- Welling, M. and Teh, Y. W. Bayesian learning via stochastic gradient langevin dynamics. In *Proceedings of the 28th international conference on machine learning (ICML-11)*, pp. 681–688, 2011.
- Zaheer, M., Kottur, S., Ravanbakhsh, S., Póczos, B., Salakhutdinov, R., and Smola, A. J. Deep sets. In *Advances in Neural Information Processing Systems*, volume 2017-December, 2017.
- Zhang, Q. and Chen, Y. Path integral sampler: a stochastic control approach for sampling. *arXiv preprint arXiv:2111.15141*, 2021.
- Štrumbelj, E., Bouchard-Côté, A., Corander, J., Gelman, A., Rue, H., Murray, L., Pesonen, H., Plummer, M., and Vehtari, A. Past, present, and future of software for bayesian inference, 2023. URL <http://hdl.handle.net/10754/694575>.

Appendix

A Related Work

In this section, we draw parallels of our work to various approaches that have been proposed to tackle the problem of either providing a good initialization for different tasks, performing implicit optimization to model predictive distributions for new tasks, or estimating the posterior through a different objective.

A.1 Variational Autoencoders

VAEs (Kingma & Welling, 2013; Rezende et al., 2014; Rezende & Mohamed, 2015; Kingma et al., 2019) are latent variable models which model observations \mathbf{x} conditioned on latent variables \mathbf{z} through the joint distribution $p_{\theta}(\mathbf{x}, \mathbf{z}) = p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})$ where $p(\mathbf{z})$ is generally chosen as $\mathcal{N}(\mathbf{0}, \mathbf{I})$. Training the model is done through VI where $q_{\varphi}(\mathbf{z})$ is obtained by explicit amortization over the data point, that is, $q_{\varphi}(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\boldsymbol{\mu}_{\varphi}(\mathbf{x}), \boldsymbol{\Sigma}_{\varphi}(\mathbf{x}))$. Training this system on a dataset \mathcal{D} is done by similarly optimizing the Evidence Lower-Bound, which boils down to the following optimization problem

$$\arg \max_{\theta, \varphi} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \mathbb{E}_{\mathbf{z} \sim q(\cdot|\mathbf{x})} \left[\log \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\varphi}(\mathbf{z}|\mathbf{x})} \right] \quad (16)$$

This objective can easily be optimized using gradient-based learning and the reparameterization trick. While typically, a diagonal Gaussian distribution is considered for q_{φ} , more complex distributions utilizing normalizing flows can also be used.

A.2 Hypernetworks

Hypernetworks are neural networks that generate weights for another neural network, used in tasks such as uncertainty quantification, zero-shot learning, etc. We refer for a comprehensive overview to Chauhan et al. (2023). Based on experiments on predicting the weights of a compact MLP (section 4), our work shows similarities with studies in this area but also has significant differences. Regarding uncertainty quantification, hypernetworks are instrumental in creating an ensemble of models by generating multiple weight vectors for the primary network. Each model within this ensemble possesses distinct parameter configurations, enabling robust estimation of uncertainty in model predictions. This feature is precious in safety-critical domains like healthcare, where confidence in predictions is essential. Multiple weight sets can be generated through techniques like dropout within hypernetworks or sampling from a noise distribution. The latter (Krueger et al., 2017) is based on a Bayesian framework where weights can be sampled using invertible network architecture, such as normalizing flows. However, while we amortize posterior inference, the weights sampled from the hypernetwork are not conditioned on information from the currently observed input data during inference time but indirectly solely on the dataset available during training, and retraining would need to be done given a new dataset. Departing from the Bayesian framework, Sun et al. (2017) have shown data-specific discriminative weight prediction, which aligns well with their specific objective of defending a convolutional neural network against adversarial attacks. Combining the ability to sample a new set of weights dataset-specifically but also handling dataset exchangeability, even in the more realistic case of missing information, our work has a distinctly different focus but also can be seen as an extension to hypernetwork research.

A.3 In-Context Learning

Amortized inference has close links to in-context learning (ICL), which has been gaining popularity, especially in natural language modeling. Various works show how in-context learning can be seen as performing implicit optimization based on the context examples, with some constructions showing exact equivalence with gradient descent in linear regression (Von Oswald et al., 2023; von Oswald et al., 2023). Other works have shown how such systems can be seen as implicitly modeling the Bayesian posterior predictive distribution (Müller et al., 2021). In a similar vein, there have been additional works aimed at directly modeling the posterior predictive distribution by providing the training data as "context" to a Transformer model and training it based on the maximum log-likelihood principle (Hollmann et al., 2022). While such approaches have been

Objective	q_φ	Model	L_2 Loss (\downarrow)					
			Gaussian		GMM			
			$2D$	$100D$	$2D-2cl$	$2D-5cl$	$5D-2cl$	$5D-5cl$
Baseline	-	Random	5.839 \pm 0.015	301.065 \pm 0.346	1.887 \pm 0.031	0.730 \pm 0.004	5.001 \pm 0.037	1.670 \pm 0.008
	-	Optimization	1.989 \pm 0.000	101.243 \pm 0.000	0.169 \pm 0.000	0.119 \pm 0.001	0.425 \pm 0.000	0.308 \pm 0.000
	-	Langevin	2.013 \pm 0.004	102.346 \pm 0.031	0.173 \pm 0.001	0.125 \pm 0.001	0.448 \pm 0.009	0.352 \pm 0.005
	-	HMC	2.018 \pm 0.008	102.413 \pm 0.028	0.174 \pm 0.001	0.135 \pm 0.001	0.479 \pm 0.007	0.449 \pm 0.002
Fwd-KL	Gaussian	GRU	2.014 \pm 0.001	102.641 \pm 0.011	0.921 \pm 0.013	0.522 \pm 0.001	2.430 \pm 0.034	1.235 \pm 0.011
		DeepSets	2.012 \pm 0.002	103.215 \pm 0.054	0.920 \pm 0.019	0.522 \pm 0.001	2.436 \pm 0.037	1.238 \pm 0.009
		Transformer	2.013 \pm 0.002	102.783 \pm 0.005	0.931 \pm 0.017	0.522 \pm 0.001	2.498 \pm 0.026	1.230 \pm 0.009
Rev-KL	Gaussian	GRU	2.012 \pm 0.001	102.509 \pm 0.008	0.183 \pm 0.002	0.132 \pm 0.002	0.471 \pm 0.010	0.413 \pm 0.019
		DeepSets	2.011 \pm 0.001	102.599 \pm 0.042	0.186 \pm 0.001	0.127 \pm 0.002	0.495 \pm 0.018	0.409 \pm 0.005
		Transformer	2.013 \pm 0.002	102.540 \pm 0.025	0.185 \pm 0.004	0.122 \pm 0.001	0.489 \pm 0.019	0.328 \pm 0.002
Fwd-KL	Flow	GRU	2.014 \pm 0.001	102.656 \pm 0.019	0.186 \pm 0.006	0.242 \pm 0.005	0.670 \pm 0.094	0.563 \pm 0.018
		DeepSets	2.014 \pm 0.001	103.340 \pm 0.029	0.185 \pm 0.006	0.237 \pm 0.008	0.648 \pm 0.082	0.583 \pm 0.028
		Transformer	2.016 \pm 0.002	102.774 \pm 0.024	0.188 \pm 0.012	0.252 \pm 0.001	0.621 \pm 0.070	0.592 \pm 0.019
Rev-KL	Flow	GRU	2.013 \pm 0.001	102.490 \pm 0.012	0.184 \pm 0.006	0.130 \pm 0.002	0.467 \pm 0.003	0.384 \pm 0.005
		DeepSets	2.011 \pm 0.001	102.674 \pm 0.046	0.188 \pm 0.005	0.131 \pm 0.002	0.519 \pm 0.008	0.405 \pm 0.005
		Transformer	2.013 \pm 0.001	102.525 \pm 0.050	0.187 \pm 0.004	0.123 \pm 0.001	0.468 \pm 0.007	0.326 \pm 0.008

Table 7: **Fixed-Dimensional.** Results for estimating the mean of a Gaussian (Gaussian) and means of a Gaussian mixture model (GMM) with the expected L_2 loss according to the posterior predictive as metric.

seeing tremendous success, they cannot be directly applied to cases where we care about and want to analyze the solution space as the solution space is only modeled implicitly, and thus, recovering it is not possible. For example, if our goal is to learn a linear regression model, an ICL model could end up learning a nonlinear model and would provide no information about the actual parameters used for prediction. As opposed to this, we obtain parameters explicitly. We thus can answer questions like the relevance of a particular feature (which corresponds to its weight in the output, and we know the weight vector explicitly). Even further, many systems grounded in physics and economics only admit a constrained solution space; for example, the movement of a human arm lies on a particular manifold, or the configuration of molecules and proteins cannot be arbitrary. Thus, performing predictions through an implicit solution space, which may violate several constraints, is not ideal. Furthermore, explicitly modeling the solution space and encoding the constraints present can be done through the prior and the parametric distribution used for modeling.

A.4 Meta Learning

Meta-learning (Hospedales et al., 2022) aims to equip models with the ability to quickly learn from different tasks or data sets to generalize to new tasks in resource-constrained domains. This attribute is precious in practical scenarios where obtaining large amounts of task-specific data is impractical or costly. A simple way of obtaining this is through nonparametric or similarity-based models like k-Nearest Neighbours, where no training is involved. Thus, new tasks can be solved quickly based on a few examples by computing a similarity metric with these examples (Koch et al., 2015; Vinyals et al., 2016; Sung et al., 2018). Another way of achieving this is through optimization-based setups, which use a nested optimization procedure. An inner step learns individual tasks from a shared initialization, whereas the outer loop computes the gradient of the whole inner process and moves the initialization in a way that allows for better generalization. Here, by relying on only a few iterations in the inner loop, the outer loop has the incentive to move the initialization to a point from which solutions to multiple tasks are reachable (Finn et al., 2017). Given the similarities between

Objective	q_φ	Model	L_2 Loss (\downarrow)		Accuracy (\uparrow)			
			Linear Regression		Linear Classification			
			$2D$	$100D$	$2D-2cl$	$2D-5cl$	$100D-2cl$	$100D-5cl$
Baseline	-	Random	4.178 \pm 0.018	202.601 \pm 0.321	50.498 \pm 0.357	19.891 \pm 0.028	50.046 \pm 0.047	20.054 \pm 0.053
	-	Optimization	0.257 \pm 0.000	25.083 \pm 0.006	96.982 \pm 0.000	93.449 \pm 0.002	70.258 \pm 0.012	41.338 \pm 0.012
	-	Langevin	0.263 \pm 0.002	23.340 \pm 0.689	95.034 \pm 0.412	88.277 \pm 0.290	65.123 \pm 0.370	32.544 \pm 0.422
	-	HMC	0.263 \pm 0.001	18.659 \pm 0.189	92.659 \pm 0.344	82.169 \pm 0.518	62.145 \pm 0.245	29.582 \pm 0.371
Fwd-KL	Gaussian	GRU	0.264 \pm 0.001	124.823 \pm 0.135	81.170 \pm 0.389	71.170 \pm 0.275	59.740 \pm 0.102	23.042 \pm 0.246
		DeepSets	0.264 \pm 0.000	123.133 \pm 1.080	81.281 \pm 0.278	70.993 \pm 0.191	50.047 \pm 0.051	20.053 \pm 0.045
		Transformer	0.264 \pm 0.000	45.856 \pm 1.331	80.960 \pm 0.285	71.484 \pm 0.437	62.954 \pm 0.062	26.789 \pm 0.110
Rev-KL	Gaussian	GRU	0.263 \pm 0.000	60.215 \pm 0.866	94.258 \pm 0.034	87.339 \pm 0.023	63.465 \pm 0.307	28.270 \pm 0.462
		DeepSets	0.263 \pm 0.000	62.837 \pm 0.617	94.285 \pm 0.116	87.342 \pm 0.021	60.867 \pm 0.265	21.339 \pm 0.085
		Transformer	0.264 \pm 0.001	28.735 \pm 0.252	94.302 \pm 0.054	87.540 \pm 0.117	68.185 \pm 0.007	32.950 \pm 0.284
Fwd-KL	Flow	GRU	0.264 \pm 0.001	119.119 \pm 0.233	96.305 \pm 0.008	88.927 \pm 0.200	59.920 \pm 0.221	23.025 \pm 0.077
		DeepSets	0.264 \pm 0.001	125.677 \pm 3.731	96.191 \pm 0.021	88.643 \pm 0.102	50.061 \pm 0.021	20.021 \pm 0.094
		Transformer	0.264 \pm 0.000	43.272 \pm 2.700	96.344 \pm 0.059	89.624 \pm 0.215	64.349 \pm 0.147	26.952 \pm 0.203
Rev-KL	Flow	GRU	0.263 \pm 0.000	61.295 \pm 1.008	95.241 \pm 0.012	88.429 \pm 0.024	64.669 \pm 0.207	28.409 \pm 1.167
		DeepSets	0.263 \pm 0.001	76.412 \pm 2.038	95.296 \pm 0.021	88.464 \pm 0.061	58.384 \pm 0.812	21.569 \pm 0.117
		Transformer	0.263 \pm 0.000	29.358 \pm 1.569	95.339 \pm 0.063	88.644 \pm 0.047	68.721 \pm 0.121	33.107 \pm 0.333

Table 8: **Fixed-Dimensional.** Results for estimating the parameters of linear regression (LR) and classification (LC) models with the expected L_2 loss and accuracy according to the posterior predictive as metrics.

meta-learning and hierarchical Bayesian inference (Grant et al., 2018), our approach can be considered as a kind of meta-learning framework; however, the line between meta-learning and Bayesian posterior inference is quite blurry as any amortized approach for the latter can be seen as a case of the former.

A.5 Neural Processes

A notable approach in meta-learning related to our research is neural processes (NP), which excel in learning scenarios with few examples. NPs (Garnelo et al., 2018a;b; Kim et al., 2019; Pakman et al., 2020; Gordon et al., 2019) can be seen as a more flexible and powerful extension of Gaussian processes that leverage a neural network-based encoder-decoder architecture for learning to model a distribution over functions that approximate a stochastic process. However, while we are interested in approximating the posterior distribution over the parameters, NPs are used to approximate the posterior predictive distribution to make predictions based on observed data. Similar to our setup, NPs rely on amortized VI for obtaining the predictive posterior. Still, instead of working with a known probabilistic model, they train the probabilistic model primarily for prediction-based tasks through approaches analogous to variational expectation maximization. Thus, they cannot provide an explicit posterior over the parameters, but they are suitable for tasks where only predictive posteriors are essential, such as those in supervised learning. NPs, in their most basic form, accomplish this by training for the objective:

$$\arg \max_{\theta, \varphi} \mathbb{E}_{\mathcal{D} \sim \chi} \mathbb{E}_{\mathbf{z} \sim q_\varphi(\cdot | \mathcal{D})} \left[\log \frac{p_\theta(\mathcal{D}, \mathbf{z})}{q_\varphi(\mathbf{z} | \mathcal{D})} \right] \quad (17)$$

where $\mathbf{z} \in \mathbb{R}^p$ is an arbitrary latent variable often uninterpretable, and the parameters of the probabilistic model θ do not get a Bayesian treatment. In particular, NPs are more suited to modeling datasets of the form $\mathcal{D} = \{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^n$, where all probabilities in Equation 17 are conditioned on the input \mathbf{x} 's, and only the predictive over \mathbf{y} 's is modeled, and p_θ is modeled as a Neural Network.

Objective	q_φ	Model	L_2 Loss (\downarrow)			
			Nonlinear Regression ReLU			
			1-layer		2-layers	
			1D	25D	1D	25D
Baseline	-	Random	65.936 \pm 0.913	831.595 \pm 8.696	1029.407 \pm 11.542	12067.691 \pm 183.598
	-	Optimization	0.360 \pm 0.001	103.967 \pm 0.110	2.370 \pm 0.015	1894.574 \pm 4.266
	-	Langevin	0.308 \pm 0.000	132.391 \pm 0.992	N/A	N/A
	-	HMC	0.374 \pm 0.002	98.061 \pm 0.730	22.314 \pm 0.814	3903.510 \pm 5.377
Fwd-KL	Gaussian	GRU	49.332 \pm 0.946	671.639 \pm 10.494	774.045 \pm 7.521	9905.246 \pm 214.545
		DeepSets	49.864 \pm 0.979	684.853 \pm 2.581	768.921 \pm 8.278	9946.090 \pm 109.933
		Transformer	49.678 \pm 0.940	680.853 \pm 5.838	747.221 \pm 12.189	9982.609 \pm 85.596
Rev-KL	Gaussian	GRU	0.426 \pm 0.004	105.976 \pm 0.586	1.066 \pm 0.069	1796.512 \pm 5.805
		DeepSets	0.426 \pm 0.004	125.853 \pm 0.791	1.394 \pm 0.108	1892.402 \pm 2.793
		Transformer	0.417 \pm 0.005	102.295 \pm 1.825	2.075 \pm 0.147	1811.440 \pm 115.435
Fwd-KL	Flow	GRU	15.781 \pm 0.210	538.962 \pm 4.269	614.925 \pm 15.494	7564.076 \pm 67.160
		DeepSets	15.051 \pm 0.120	548.535 \pm 3.288	622.461 \pm 7.043	7618.364 \pm 115.946
		Transformer	16.109 \pm 0.307	539.338 \pm 4.336	597.718 \pm 8.358	7635.052 \pm 109.037
Rev-KL	Flow	GRU	0.405 \pm 0.010	106.001 \pm 0.420	0.988 \pm 0.045	1814.649 \pm 8.327
		DeepSets	0.395 \pm 0.004	128.169 \pm 1.451	1.215 \pm 0.028	1886.698 \pm 7.294
		Transformer	0.387 \pm 0.004	102.610 \pm 0.863	2.549 \pm 0.058	1791.741 \pm 49.585

Table 9: **Fixed-Dimensional.** Results for estimating the parameters of nonlinear regression models with ReLU activation function, with the expected L_2 loss according to the posterior predictive as metric.

These approaches can be seen as quite related to ICL, where the exchangeable architecture backbone is switched from DeepSets to Transformers. Similar to ICL, they do not provide control over the solution space as they aim to model either the posterior predictive or an arbitrary latent space. While this leads to good predictive performance on various tasks, they cannot be freely applied to problems that pose certain constraints on the underlying probabilistic model. In such cases, estimating the actual parameters is important to enforce constraints in the parameter space as well as for interpretability, which we already discussed in the ICL section.

A.6 Simulation-Based Inference

In the case of simulation-based inference (Cranmer et al., 2020), when the likelihood $p(\mathbf{x}|\boldsymbol{\theta})$ is intractable, BayesFlow (Radev et al., 2020) and similar methods (Lorch et al., 2022) provide a solution framework to amortize Bayesian inference of parameters in complex models. Starting from the forward KL divergence between the true and approximate posteriors, the resulting objective is to optimize for parameters of the approximate posterior distribution that maximize the posterior probability of data-generating parameters $\boldsymbol{\theta}$ given observed data \mathcal{D} for all $\boldsymbol{\theta}$ and \mathcal{D} . Density estimation of the approximate posterior can then be done using the change-of-variables formula and a conditional invertible neural network that parameterizes the approximate posterior distribution.

$$\arg \min_{\varphi} \mathbb{KL}[p(\boldsymbol{\theta}|\mathcal{D})||q_{\varphi}(\boldsymbol{\theta}|\mathcal{D})] = \arg \min_{\varphi=\{\nu,\psi\}} \mathbb{E}_{(\boldsymbol{\theta},\mathcal{D})\sim p(\boldsymbol{\theta},\mathcal{D})} [-\log p_{\mathbf{z}}(f_{\nu}(\boldsymbol{\theta}; h_{\psi}(\mathcal{D}))) - \log |\det J_{f_{\nu}}|] \quad (18)$$

Since their goal is to learn a global estimator for the probabilistic mapping from \mathcal{D} to data generating $\boldsymbol{\theta}$, the information about the observed dataset is encoded in the output of a summary network h_{ψ} . It is used as conditional input to the normalizing flow f_{ν} . Although the likelihood function does not need to be

Objective	q_φ	Model	L_2 Loss (\downarrow)			
			Nonlinear Regression TanH			
			1-layer		2-layers	
			1D	25D	1D	25D
Baseline	-	Random	31.448 \pm 0.186	52.644 \pm 0.173	52.735 \pm 1.122	52.583 \pm 0.132
	-	Optimization	0.366 \pm 0.001	13.352 \pm 0.005	0.651 \pm 0.002	30.176 \pm 0.056
	-	Langevin	0.296 \pm 0.003	17.221 \pm 0.130	0.363 \pm 0.003	28.528 \pm 0.115
	-	HMC	0.398 \pm 0.003	12.607 \pm 0.192	0.733 \pm 0.021	23.571 \pm 0.346
Fwd-KL	Gaussian	GRU	31.391 \pm 0.161	52.008 \pm 0.282	52.725 \pm 1.149	51.989 \pm 0.139
		DeepSets	31.421 \pm 0.074	52.137 \pm 0.215	52.850 \pm 1.192	51.904 \pm 0.334
		Transformer	31.350 \pm 0.219	52.945 \pm 0.430	52.693 \pm 1.188	52.364 \pm 0.164
Rev-KL	Gaussian	GRU	0.415 \pm 0.003	15.874 \pm 6.958	0.951 \pm 0.047	25.907 \pm 0.012
		DeepSets	0.405 \pm 0.004	25.333 \pm 0.010	0.912 \pm 0.013	25.877 \pm 0.002
		Transformer	0.412 \pm 0.013	11.784 \pm 0.949	0.847 \pm 0.010	20.405 \pm 3.874
Fwd-KL	Flow	GRU	12.415 \pm 0.800	52.039 \pm 0.065	52.695 \pm 0.611	52.576 \pm 0.225
		DeepSets	31.790 \pm 0.163	51.933 \pm 0.115	52.903 \pm 0.625	52.643 \pm 0.239
		Transformer	10.392 \pm 0.195	52.470 \pm 0.364	52.385 \pm 0.689	52.646 \pm 0.622
Rev-KL	Flow	GRU	0.386 \pm 0.005	11.401 \pm 0.041	0.736 \pm 0.009	25.892 \pm 0.010
		DeepSets	0.374 \pm 0.005	25.685 \pm 0.004	0.686 \pm 0.019	25.885 \pm 0.007
		Transformer	0.376 \pm 0.002	10.486 \pm 0.040	0.724 \pm 0.026	25.885 \pm 0.011

Table 10: **Fixed-Dimensional.** Results for estimating the parameters of nonlinear regression models with TanH activation function, with the expected L_2 loss according to the posterior predictive as metric.

known, the method requires access to paired observations $(\mathbf{x}, \boldsymbol{\theta})$ for training, which is sometimes unavailable. This approach is equivalent to the *Forward KL* setup in our experiments when trained with DeepSets and Normalizing Flows. Current research has also leveraged score-based generative models for SBI which can condition on a dataset by learning a score model conditional only on single observations (Geffner et al., 2023).

A.7 Amortization in Gaussian Processes

Gaussian Processes (GPs) define a class of probabilistic models that do enjoy tractable likelihood. However, inference in such systems is slow and sensitive to the choice of kernel function that defines the covariance matrix. Similar to meta learning and neural processes, current research also focuses on estimating the kernel function in GPs by leveraging permutation invariant architectures like transformers (Liu et al., 2020; Simpson et al., 2021; Bitzer et al., 2023). Additionally, often these approaches amortize based on point estimates and are leveraged when considering GPs for regression problems, and it is not straightforward to extend them to classification or unsupervised learning. In contrast, our approach is more general and can work for all problems that define a differentiable likelihood function. Additionally, our approach also approximates the Bayesian posterior distribution over the parameters of interest, as opposed to point estimates.

A.8 Mode Collapse in Variational Inference

Reverse KL based methods have been widely known to suffer from mode collapse due to the nature of the optimization objective (Bishop & Nasrabadi, 2006), which implies that even if the approximate distribution possesses the ability to represent multiple modes, optimization is often sub-optimal and the distribution ends up covering only a small handful of them. Improving normalizing flow based methods with repulsive terms or through the lens of natural gradient optimization procedure for a mixture approximate distribution (Arenz

Objective	q_φ	Model	Accuracy (\uparrow)			
			Nonlinear Classification ReLU - 2 class			
			1-layer		2-layers	
			2D	25D	2D	25D
Baseline	-	Random	50.306 \pm 0.590	50.008 \pm 0.326	50.394 \pm 0.190	49.846 \pm 0.635
	-	Optimization	96.879 \pm 0.028	77.896 \pm 0.023	96.770 \pm 0.062	82.073 \pm 0.200
	-	Langevin	95.971 \pm 0.313	73.165 \pm 0.282	96.645 \pm 0.101	76.541 \pm 0.139
	-	HMC	91.763 \pm 0.163	70.395 \pm 0.110	91.797 \pm 0.048	76.445 \pm 0.372
Fwd-KL	Gaussian	GRU	59.518 \pm 0.355	56.858 \pm 0.319	60.962 \pm 0.599	60.063 \pm 0.695
		DeepSets	59.383 \pm 0.244	56.806 \pm 0.204	61.090 \pm 0.690	59.933 \pm 0.599
		Transformer	59.588 \pm 0.389	57.089 \pm 0.376	61.151 \pm 0.560	60.041 \pm 0.680
Rev-KL	Gaussian	GRU	92.384 \pm 0.195	72.455 \pm 0.032	86.157 \pm 0.066	69.966 \pm 0.333
		DeepSets	92.488 \pm 0.133	59.806 \pm 0.315	86.275 \pm 0.733	69.550 \pm 0.371
		Transformer	92.627 \pm 0.377	75.178 \pm 0.142	86.351 \pm 0.217	69.812 \pm 0.527
Fwd-KL	Flow	GRU	76.931 \pm 0.266	58.338 \pm 0.026	62.981 \pm 0.452	63.199 \pm 0.202
		DeepSets	72.313 \pm 1.829	58.113 \pm 0.126	62.438 \pm 0.277	62.884 \pm 0.201
		Transformer	77.296 \pm 0.201	58.344 \pm 0.148	63.753 \pm 0.155	63.590 \pm 0.390
Rev-KL	Flow	GRU	93.392 \pm 0.073	72.002 \pm 0.506	85.712 \pm 0.665	71.419 \pm 0.200
		DeepSets	93.312 \pm 0.197	60.943 \pm 0.184	85.960 \pm 0.896	71.288 \pm 0.211
		Transformer	93.578 \pm 0.093	74.956 \pm 0.546	87.138 \pm 0.438	71.525 \pm 0.159

Table 11: **Fixed-Dimensional.** Results for estimating the parameters of nonlinear classification models with ReLU activation function and two classes, with the expected accuracy according to the posterior predictive as metric.

et al., 2022; Lin et al., 2020) is an important topic of research, and we believe it would be quite an important future work to experimentally validate if they help in learning multi-modality in amortized posterior inference problems that are studied in this work.

B Architectures respecting Exchangeability

In this section, we highlight how DeepSets and Transformer models satisfy the dataset exchangeability criteria, which is essential in modeling the posterior distribution over the parameters of any probabilistic model relying on *iid* data.

B.1 DeepSets

DeepSets (Zaheer et al., 2017) operate on arbitrary sets $\mathcal{X} = \{x_1, \dots, x_N\} \subset \mathbb{R}^d$ of fixed dimensionality d by first mapping each individual element $x_i \in \mathcal{X}$ to some high-dimensional space using a nonlinear transform, which is parameterized as a multi-layered neural network with parameters φ_1

$$z_i = f_{\varphi_1}(x_i) \quad (19)$$

After having obtained this high-dimensional embedding of each element of the set, it applies an aggregation function $a(\cdot)$, which is a permutation invariant function that maps a set of elements $\mathcal{Z} = \{z_1, \dots, z_N\} \in \mathbb{R}^z$ to an element $\mathbf{h} \in \mathbb{R}^z$,

$$\mathbf{h} = a(\mathcal{Z}) \quad (20)$$

Objective	q_φ	Model	Accuracy (\uparrow)			
			Nonlinear Classification ReLU - 5 class			
			1-layer		2-layers	
			2D	25D	2D	25D
Baseline	-	Random	19.952 \pm 0.119	20.026 \pm 0.062	19.832 \pm 0.202	19.985 \pm 0.091
	-	Optimization	94.369 \pm 0.022	60.589 \pm 0.064	93.664 \pm 0.025	60.824 \pm 0.023
	-	Langevin	91.286 \pm 0.149	50.845 \pm 0.366	92.449 \pm 0.068	52.216 \pm 0.170
	-	HMC	81.387 \pm 0.742	47.408 \pm 0.548	81.098 \pm 0.246	52.854 \pm 0.190
Fwd-KL	Gaussian	GRU	32.597 \pm 0.169	29.932 \pm 0.142	31.061 \pm 0.089	30.187 \pm 0.070
		DeepSets	32.481 \pm 0.090	29.909 \pm 0.218	30.923 \pm 0.205	19.983 \pm 0.089
		Transformer	32.977 \pm 0.114	30.064 \pm 0.134	31.478 \pm 0.079	30.307 \pm 0.121
Rev-KL	Gaussian	GRU	83.460 \pm 0.269	36.688 \pm 0.189	71.710 \pm 0.375	28.234 \pm 0.148
		DeepSets	83.734 \pm 0.313	34.078 \pm 0.066	67.519 \pm 1.735	47.585 \pm 0.143
		Transformer	84.645 \pm 0.515	36.850 \pm 0.138	74.390 \pm 0.280	27.808 \pm 0.232
Fwd-KL	Flow	GRU	43.850 \pm 0.412	31.926 \pm 0.096	33.771 \pm 0.400	31.871 \pm 0.228
		DeepSets	43.620 \pm 0.250	31.829 \pm 0.126	32.919 \pm 0.126	19.950 \pm 0.192
		Transformer	44.057 \pm 0.122	31.789 \pm 0.094	34.079 \pm 0.224	32.392 \pm 0.212
Rev-KL	Flow	GRU	84.459 \pm 0.194	37.343 \pm 0.134	64.309 \pm 0.273	48.459 \pm 0.542
		DeepSets	84.775 \pm 0.129	34.917 \pm 0.069	67.937 \pm 4.341	48.563 \pm 0.199
		Transformer	85.857 \pm 0.096	46.968 \pm 0.094	75.005 \pm 0.605	48.889 \pm 0.380

Table 12: **Fixed-Dimensional.** Results for estimating the parameters of nonlinear classification models with ReLU activation function and five classes, with the expected accuracy according to the posterior predictive as metric.

Thus, the outcome does not change under permutations of \mathcal{Z} . Finally, another nonlinear transform, parameterized by a multi-layered neural network with parameters φ_2 , is applied to the outcome \mathbf{h} to provide the final output.

$$\mathbf{o} = g_{\varphi_2}(\mathbf{h}) \quad (21)$$

For our experiments, we then use the vector \mathbf{o} to predict the parameters of a parametric family of distributions (e.g., Gaussian or Flows) using an additional nonlinear neural network. As an example, for the Gaussian case, we consider the distribution $\mathcal{N}(\cdot|\boldsymbol{\mu}, \boldsymbol{\Sigma})$, where

$$\boldsymbol{\mu} := \boldsymbol{\mu}_{\varphi_3}(\mathbf{o}) \quad \text{and} \quad \boldsymbol{\Sigma} := \boldsymbol{\Sigma}_{\varphi_4}(\mathbf{o}) \quad (22)$$

which makes $\boldsymbol{\mu}$ implicitly a function of the original input set \mathcal{X} . To understand why the posterior distribution modeled in this fashion does not change when the inputs are permuted, let us assume that Π is a permutation over the elements of \mathcal{X} . If we look at one of the parameters of the posterior distribution, e.g., $\boldsymbol{\mu}$, we can see that

$$\boldsymbol{\mu}(\Pi\mathcal{X}) = \boldsymbol{\mu}_{\varphi_3}(g_{\varphi_2}(a(\{f_{\varphi_1}(\mathbf{x}_{\Pi(i)})\}_{i=1}^N))) \quad (23)$$

$$= \boldsymbol{\mu}_{\varphi_3}(g_{\varphi_2}(a(\{f_{\varphi_1}(\mathbf{x}_i)\}_{i=1}^N))) \quad (24)$$

$$= \boldsymbol{\mu}(\mathcal{X}) \quad (25)$$

which simply follows from the fact that $a(\cdot)$ is a permutation invariant operation, e.g., sum or mean. We can also provide similar reasoning for the other parameters (e.g., $\boldsymbol{\Sigma}$). This shows that DeepSets can be used to model the posterior distribution over parameters of interest as it respects the exchangeability criteria (*iid* observations) assumptions in the data through its permutation invariant structure.

Objective	q_φ	Model	Accuracy (\uparrow)			
			Nonlinear Classification TanH - 2 class			
			1-layer		2-layers	
			2D	25D	2D	25D
Baseline	-	Random	50.147 \pm 0.603	49.963 \pm 0.083	49.942 \pm 0.452	49.978 \pm 0.113
	-	Optimization	96.552 \pm 0.005	75.228 \pm 0.029	94.130 \pm 0.018	69.052 \pm 0.029
	-	Langevin	94.778 \pm 0.210	68.787 \pm 0.258	92.417 \pm 0.149	62.722 \pm 0.180
	-	HMC	91.674 \pm 0.200	67.415 \pm 0.667	88.218 \pm 0.202	62.323 \pm 0.262
Fwd-KL	Gaussian	GRU	50.151 \pm 0.616	49.959 \pm 0.087	49.942 \pm 0.448	49.987 \pm 0.112
		DeepSets	50.147 \pm 0.601	49.961 \pm 0.083	49.939 \pm 0.451	49.977 \pm 0.110
		Transformer	50.146 \pm 0.622	49.959 \pm 0.082	49.948 \pm 0.448	49.970 \pm 0.104
Rev-KL	Gaussian	GRU	89.813 \pm 0.181	49.969 \pm 0.085	49.956 \pm 0.442	49.974 \pm 0.088
		DeepSets	89.558 \pm 0.264	49.970 \pm 0.100	49.961 \pm 0.436	49.986 \pm 0.094
		Transformer	89.879 \pm 0.387	49.978 \pm 0.075	79.102 \pm 0.155	49.987 \pm 0.089
Fwd-KL	Flow	GRU	50.066 \pm 0.488	49.835 \pm 0.240	49.836 \pm 0.328	49.911 \pm 0.206
		DeepSets	50.078 \pm 0.468	49.848 \pm 0.242	49.834 \pm 0.333	49.910 \pm 0.205
		Transformer	50.067 \pm 0.483	49.827 \pm 0.235	49.813 \pm 0.331	49.905 \pm 0.210
Rev-KL	Flow	GRU	90.396 \pm 0.126	49.910 \pm 0.169	50.018 \pm 0.171	50.001 \pm 0.076
		DeepSets	90.247 \pm 0.067	49.901 \pm 0.124	50.132 \pm 0.162	49.942 \pm 0.110
		Transformer	90.416 \pm 0.311	49.904 \pm 0.133	81.729 \pm 0.070	49.945 \pm 0.203

Table 13: **Fixed-Dimensional.** Results for estimating the parameters of nonlinear classification models with TanH activation function and two classes, with the expected accuracy according to the posterior predictive as metric.

B.2 Transformers

Similarly, we can look at Transformers (Vaswani et al., 2017) as candidates for respecting the exchangeability conditions in the data. In particular, we consider transformer systems without positional encodings and consider an additional [CLS] token, denoted by $\mathbf{c} \in \mathbb{R}^d$, to drive the prediction. If we look at the application of a layer of transformer model, it can be broken down into two components.

Multi-Head Attention. Given a query vector obtained from \mathbf{c} and keys and values coming from our input set $\mathcal{X} \subset \mathbb{R}^d$, we can model the update of the context \mathbf{c} as

$$\hat{\mathbf{c}}(\mathcal{X}) = \text{Softmax}(\mathbf{c}^T \mathbf{W}_Q \mathbf{W}_K^T \mathbf{X}^T) \mathbf{X} \mathbf{W}_V \quad (26)$$

where $\mathbf{W}_Q \in \mathbb{R}^{d \times k}$, $\mathbf{W}_K \in \mathbb{R}^{d \times k}$, $\mathbf{W}_V \in \mathbb{R}^{d \times k}$ and $\mathbf{X} \in \mathbb{R}^{N \times d}$ denotes a certain ordering of the elements in \mathcal{X} . Further, $\hat{\mathbf{c}}$ is the updated vector after attention, and Softmax is over the rows of \mathbf{X} . Here, we see that if we were to apply a permutation to the elements in \mathbf{X} , the outcome would remain the same. In particular

$$\hat{\mathbf{c}}(\Pi \mathbf{X}) = \text{Softmax}(\mathbf{c}^T \mathbf{W}_Q \mathbf{W}_K^T \mathbf{X}^T \Pi^T) \Pi \mathbf{X} \mathbf{W}_V \quad (27)$$

$$= \text{Softmax}(\mathbf{c}^T \mathbf{W}_Q \mathbf{W}_K^T \mathbf{X}^T) \Pi^T \Pi \mathbf{X} \mathbf{W}_V \quad (28)$$

$$= \text{Softmax}(\mathbf{c}^T \mathbf{W}_Q \mathbf{W}_K^T \mathbf{X}^T) \mathbf{X} \mathbf{W}_V \quad (29)$$

$$= \hat{\mathbf{c}}(\mathbf{X}) \quad (30)$$

which follows because Softmax is an equivariant function, i.e., applying Softmax on a permutation of columns is equivalent to applying Softmax first and then permuting the columns correspondingly. Thus, we see that the update to the [CLS] token \mathbf{c} is permutation invariant. This output is then used independently as input to

Objective	q_φ	Model	Accuracy (\uparrow)			
			Nonlinear Classification TanH - 5 class			
			1-layer		2-layers	
			2D	25D	2D	25D
Baseline	-	Random	19.745 \pm 0.269	20.037 \pm 0.075	20.214 \pm 0.043	19.864 \pm 0.089
	-	Optimization	92.919 \pm 0.011	49.972 \pm 0.070	88.156 \pm 0.014	39.410 \pm 0.043
	-	Langevin	88.635 \pm 0.374	39.758 \pm 0.212	84.050 \pm 0.138	31.472 \pm 0.010
	-	HMC	81.057 \pm 0.277	35.378 \pm 0.130	75.305 \pm 0.134	29.730 \pm 0.670
Fwd-KL	Gaussian	GRU	19.960 \pm 0.271	20.235 \pm 0.083	20.452 \pm 0.059	20.028 \pm 0.077
		DeepSets	19.807 \pm 0.209	20.040 \pm 0.072	20.210 \pm 0.043	19.861 \pm 0.094
		Transformer	19.977 \pm 0.273	20.241 \pm 0.068	20.453 \pm 0.062	20.029 \pm 0.082
Rev-KL	Gaussian	GRU	77.711 \pm 0.014	20.026 \pm 0.079	20.213 \pm 0.043	19.877 \pm 0.092
		DeepSets	76.414 \pm 0.378	20.038 \pm 0.060	20.216 \pm 0.027	19.887 \pm 0.089
		Transformer	79.163 \pm 0.183	20.026 \pm 0.093	51.408 \pm 0.484	19.872 \pm 0.083
Fwd-KL	Flow	GRU	32.900 \pm 0.115	20.209 \pm 0.048	20.105 \pm 0.385	20.040 \pm 0.037
		DeepSets	20.137 \pm 0.070	20.000 \pm 0.048	19.887 \pm 0.392	19.895 \pm 0.025
		Transformer	30.135 \pm 2.459	20.224 \pm 0.048	20.104 \pm 0.399	20.030 \pm 0.039
Rev-KL	Flow	GRU	79.329 \pm 0.320	20.074 \pm 0.033	19.904 \pm 0.180	19.864 \pm 0.079
		DeepSets	20.071 \pm 0.302	19.999 \pm 0.045	19.872 \pm 0.259	19.906 \pm 0.051
		Transformer	80.064 \pm 0.159	20.002 \pm 0.032	19.777 \pm 0.210	19.911 \pm 0.082

Table 14: **Fixed-Dimensional**. Results for estimating the parameters of nonlinear classification models with TanH activation function and five classes, with the expected accuracy according to the posterior predictive as metric.

a multi-layered neural network with residual connections, and the entire process is repeated multiple times without weight sharing to simulate multiple layers. Since all the individual parts are permutation invariant w.r.t permutations on \mathcal{X} , the entire setup ends up being permutation invariant. Obtaining the parameters of a parametric family of distribution for posterior estimation then follows the same recipe as DeepSets, with \mathbf{o} replaced by \mathbf{c} .

C Probabilistic Models

This section details the various candidate probabilistic models used in our experiments for amortized computation of Bayesian posteriors over the parameters. Here, we explain the parameters associated with the probabilistic model over which we want to estimate the posterior and the likelihood and prior that we use for experimentation.

Mean of Gaussian (GM): As a proof of concept, we consider the simple setup of estimating the posterior distribution over the mean of a Gaussian distribution $p(\boldsymbol{\mu}|\mathcal{D})$ given some observed data. In this case, prior and likelihood defining the probabilistic model $p(\mathbf{x}, \boldsymbol{\theta})$ (with $\boldsymbol{\theta}$ being the mean $\boldsymbol{\mu}$) are given by:

$$p(\boldsymbol{\mu}) = \mathcal{N}(\boldsymbol{\mu}|\mathbf{0}, \mathbf{I}) \tag{31}$$

$$p(\mathbf{x}|\boldsymbol{\mu}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) \tag{32}$$

and $\boldsymbol{\Sigma}$ is known beforehand and defined as a unit variance matrix.

Objective	q_φ	Model	L_2 Loss (\downarrow)					
			Gaussian		GMM			
			$2D$	$100D$	$2D-2cl$	$2D-5cl$	$5D-2cl$	$5D-5cl$
Baseline	-	Random	6.297 \pm 0.017	298.238 \pm 0.228	2.078 \pm 0.134	0.626 \pm 0.037	4.659 \pm 0.034	1.632 \pm 0.004
	-	Optimization	2.020 \pm 0.000	100.885 \pm 0.000	0.175 \pm 0.002	0.121 \pm 0.002	0.427 \pm 0.000	0.323 \pm 0.002
	-	Langevin	2.036 \pm 0.004	101.917 \pm 0.042	0.178 \pm 0.002	0.123 \pm 0.002	0.440 \pm 0.002	0.340 \pm 0.006
	-	HMC	2.044 \pm 0.008	102.015 \pm 0.009	0.189 \pm 0.013	0.132 \pm 0.004	0.462 \pm 0.013	0.423 \pm 0.009
Fwd-KL	Gaussian	GRU	2.300 \pm 0.105	133.224 \pm 0.579	1.119 \pm 0.150	0.473 \pm 0.012	2.360 \pm 0.017	1.208 \pm 0.002
		DeepSets	2.216 \pm 0.017	129.695 \pm 0.737	1.113 \pm 0.150	0.477 \pm 0.013	2.352 \pm 0.018	1.210 \pm 0.003
		Transformer	2.352 \pm 0.013	108.977 \pm 0.100	1.134 \pm 0.153	0.476 \pm 0.013	2.399 \pm 0.019	1.208 \pm 0.004
Rev-KL	Gaussian	GRU	2.047 \pm 0.002	105.141 \pm 0.102	0.187 \pm 0.005	0.147 \pm 0.006	0.462 \pm 0.011	0.398 \pm 0.019
		DeepSets	2.049 \pm 0.003	105.062 \pm 0.212	0.194 \pm 0.004	0.145 \pm 0.003	0.481 \pm 0.021	0.387 \pm 0.002
		Transformer	2.057 \pm 0.004	104.709 \pm 0.122	0.195 \pm 0.004	0.140 \pm 0.003	0.468 \pm 0.006	0.335 \pm 0.016
Fwd-KL	Flow	GRU	2.358 \pm 0.005	125.835 \pm 1.983	0.287 \pm 0.032	0.212 \pm 0.002	0.596 \pm 0.068	0.513 \pm 0.015
		DeepSets	2.053 \pm 0.003	133.229 \pm 1.933	0.271 \pm 0.050	0.202 \pm 0.006	0.584 \pm 0.030	0.517 \pm 0.016
		Transformer	2.060 \pm 0.005	108.484 \pm 0.164	0.344 \pm 0.054	0.221 \pm 0.006	0.591 \pm 0.080	0.533 \pm 0.010
Rev-KL	Flow	GRU	2.050 \pm 0.004	105.187 \pm 0.030	0.199 \pm 0.014	0.142 \pm 0.004	0.466 \pm 0.007	0.373 \pm 0.004
		DeepSets	2.054 \pm 0.003	105.095 \pm 0.064	0.202 \pm 0.007	0.146 \pm 0.003	0.494 \pm 0.013	0.379 \pm 0.003
		Transformer	2.049 \pm 0.003	104.914 \pm 0.113	0.193 \pm 0.004	0.138 \pm 0.002	0.460 \pm 0.003	0.327 \pm 0.006

Table 15: **Variable-Dimensional.** Results for estimating the mean of a Gaussian (Gaussian) and means of a Gaussian mixture model (GMM) with the expected L_2 loss according to the posterior predictive as metric.

Linear Regression (LR): We then look at the problem of estimating the posterior over the weight vector for Bayesian linear regression given a dataset $p(\mathbf{w}, b|\mathcal{D})$, where the underlying model $p(\mathcal{D}, \boldsymbol{\theta})$ is given by:

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \mathbf{I}) \quad (33)$$

$$p(b) = \mathcal{N}(b|0, 1) \quad (34)$$

$$p(y|\mathbf{x}, \mathbf{w}, b) = \mathcal{N}(y|\mathbf{w}^T \mathbf{x} + b, \sigma^2), \quad (35)$$

and with $\sigma^2 = 0.25$ known beforehand. Inputs \mathbf{x} are generated from $p(\mathbf{x}) = \mathcal{N}(\mathbf{0}, I)$.

Linear Classification (LC): We now consider a setting where the true posterior cannot be obtained analytically as the likelihood and prior are not conjugate. In this case, we consider the underlying probabilistic model by:

$$p(\mathbf{W}) = \mathcal{N}(\mathbf{W}|\mathbf{0}, \mathbf{I}) \quad (36)$$

$$p(y|\mathbf{x}, \mathbf{W}) = \text{Categorical}\left(y \left| \frac{1}{\tau} \mathbf{W} \mathbf{x} \right. \right), \quad (37)$$

where τ is the known temperature term which is kept as 0.1 to ensure peaky distributions, and \mathbf{x} is being generated from $p(\mathbf{x}) = \mathcal{N}(\mathbf{0}, I)$.

Nonlinear Regression (NLR): Next, we tackle the more complex problem where the posterior distribution is multi-modal and obtaining multiple modes or even a single good one is challenging. For this, we consider the model as a Bayesian Neural Network (BNN) for regression with fixed hyper-parameters like the number of layers, dimensionality of the hidden layer, etc. Let the BNN denote the function $f_{\boldsymbol{\theta}}$ where $\boldsymbol{\theta}$ are the network parameters such that the estimation problem is to approximate $p(\boldsymbol{\theta}|\mathcal{D})$. Then, for regression, we specify the

Objective	q_φ	Model	L_2 Loss (\downarrow)		Accuracy (\uparrow)			
			Linear Regression		Linear Classification			
			$2D$	$100D$	$2D-2cl$	$2D-5cl$	$100D-2cl$	$100D-5cl$
Baseline	-	Random	4.272 \pm 0.068	200.836 \pm 0.609	50.125 \pm 0.264	20.078 \pm 0.065	50.005 \pm 0.061	20.033 \pm 0.082
	-	Optimization	0.258 \pm 0.000	20.127 \pm 0.003	97.301 \pm 0.000	91.752 \pm 0.000	71.231 \pm 0.010	42.345 \pm 0.001
	-	Langevin	0.263 \pm 0.002	21.781 \pm 0.953	95.441 \pm 0.209	86.445 \pm 0.496	65.469 \pm 0.513	32.668 \pm 0.145
	-	HMC	0.263 \pm 0.000	17.774 \pm 0.120	92.961 \pm 0.228	78.793 \pm 0.314	62.602 \pm 0.171	30.055 \pm 0.506
Fwd-KL	Gaussian	GRU	0.271 \pm 0.004	139.396 \pm 1.012	79.467 \pm 0.711	65.124 \pm 0.861	57.872 \pm 0.157	22.677 \pm 0.081
		DeepSets	0.269 \pm 0.001	149.784 \pm 0.766	80.323 \pm 0.429	20.078 \pm 0.059	50.767 \pm 0.058	20.035 \pm 0.081
		Transformer	0.279 \pm 0.001	64.282 \pm 3.711	79.901 \pm 0.271	60.984 \pm 1.590	62.382 \pm 0.029	26.997 \pm 0.098
Rev-KL	Gaussian	GRU	0.291 \pm 0.013	62.624 \pm 0.123	93.367 \pm 0.289	82.020 \pm 0.127	63.411 \pm 0.248	28.655 \pm 0.149
		DeepSets	0.279 \pm 0.004	64.064 \pm 0.221	93.977 \pm 0.093	83.832 \pm 0.106	61.305 \pm 0.114	27.877 \pm 0.265
		Transformer	0.271 \pm 0.007	31.984 \pm 0.482	94.336 \pm 0.210	82.976 \pm 0.074	67.676 \pm 0.078	33.125 \pm 0.051
Fwd-KL	Flow	GRU	0.273 \pm 0.001	138.284 \pm 1.030	92.078 \pm 0.190	75.151 \pm 0.274	57.982 \pm 0.138	22.430 \pm 0.208
		DeepSets	0.270 \pm 0.002	153.207 \pm 0.814	85.950 \pm 6.332	20.059 \pm 0.233	50.494 \pm 0.055	19.976 \pm 0.095
		Transformer	0.276 \pm 0.004	63.102 \pm 1.963	94.494 \pm 0.368	74.876 \pm 0.857	63.559 \pm 0.072	27.098 \pm 0.147
Rev-KL	Flow	GRU	0.276 \pm 0.008	71.260 \pm 1.265	94.292 \pm 0.175	83.622 \pm 0.057	63.391 \pm 0.133	27.340 \pm 0.243
		DeepSets	0.274 \pm 0.003	76.772 \pm 1.836	94.570 \pm 0.178	85.059 \pm 0.127	59.116 \pm 0.491	22.810 \pm 0.245
		Transformer	0.279 \pm 0.013	33.056 \pm 0.321	94.793 \pm 0.135	84.929 \pm 0.027	68.124 \pm 0.214	33.251 \pm 0.130

Table 16: **Variable-Dimensional.** Results for estimating the parameters of linear regression (LR) and classification (LC) models with the expected L_2 loss and accuracy according to the posterior predictive as metrics.

probabilistic model using:

$$p(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta}|\mathbf{0}, \mathbf{I}) \quad (38)$$

$$p(y|\mathbf{x}, \boldsymbol{\theta}) = \mathcal{N}(y|f_{\boldsymbol{\theta}}(\mathbf{x}), \sigma^2), \quad (39)$$

where $\sigma^2 = 0.25$ is a known quantity and \mathbf{x} being generated from $p(\mathbf{x}) = \mathcal{N}(\mathbf{0}, I)$.

Nonlinear Classification (NLC): Like in Nonlinear Regression, we consider BNNs with fixed hyperparameters for classification problems with the same estimation task of approximating $p(\boldsymbol{\theta}|\mathcal{D})$. In this formulation, we consider the probabilistic model as:

$$p(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta}|\mathbf{0}, \mathbf{I}) \quad (40)$$

$$p(y|\mathbf{x}, \boldsymbol{\theta}) = \text{Categorical}\left(y \left| \frac{1}{\tau} f_{\boldsymbol{\theta}}(\mathbf{x})\right.\right) \quad (41)$$

where τ is the known temperature term which is kept as 0.1 to ensure peaky distributions, and \mathbf{x} is being generated from $p(\mathbf{x}) = \mathcal{N}(\mathbf{0}, I)$.

Gaussian Mixture Model (GMM): While we have mostly looked at predictive problems, where the task is to model some predictive variable y conditioned on some input \mathbf{x} , we now look at a well-known probabilistic model for unsupervised learning, Gaussian Mixture Model (GMM), primarily used to cluster data. Consider a K -cluster GMM with:

$$p(\boldsymbol{\mu}_k) = \mathcal{N}(\boldsymbol{\mu}_k|\mathbf{0}, \mathbf{I}) \quad (42)$$

$$p(\mathbf{x}|\boldsymbol{\mu}_{1:K}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k). \quad (43)$$

Objective	q_φ	Model	L_2 Loss (\downarrow)			
			1-layer		2-layers	
			1D	50D	1D	50D
Baseline	-	Random	73.006 \pm 0.171	1704.335 \pm 9.330	998.890 \pm 16.317	27799.887 \pm 165.653
	-	Optimization	0.359 \pm 0.002	309.162 \pm 0.204	3.078 \pm 0.057	4894.141 \pm 4.290
	-	Langevin	0.308 \pm 0.002	N/A	N/A	N/A
	-	HMC	0.381 \pm 0.005	303.857 \pm 2.487	7.999 \pm 0.595	12905.846 \pm 9.903
Fwd-KL	Gaussian	GRU	51.448 \pm 0.029	1346.462 \pm 6.833	754.388 \pm 32.650	21121.735 \pm 112.458
		DeepSets	51.901 \pm 1.542	1357.462 \pm 5.348	767.781 \pm 18.262	21110.264 \pm 72.613
		Transformer	50.813 \pm 0.532	1319.868 \pm 12.165	753.177 \pm 23.023	21037.666 \pm 116.929
Rev-KL	Gaussian	GRU	2.308 \pm 0.126	316.344 \pm 6.189	23.520 \pm 3.816	4673.633 \pm 98.443
		DeepSets	0.977 \pm 0.118	451.942 \pm 2.785	8.034 \pm 0.375	6127.317 \pm 190.224
		Transformer	0.815 \pm 0.024	278.282 \pm 1.073	8.300 \pm 1.673	4744.375 \pm 24.609
Fwd-KL	Flow	GRU	38.407 \pm 0.364	1097.410 \pm 9.498	664.409 \pm 9.590	18372.905 \pm 62.019
		DeepSets	43.305 \pm 2.063	1119.990 \pm 5.461	745.412 \pm 27.599	20719.466 \pm 627.997
		Transformer	39.701 \pm 0.519	1073.256 \pm 1.504	619.631 \pm 22.299	16700.463 \pm 350.596
Rev-KL	Flow	GRU	2.305 \pm 0.406	302.918 \pm 5.644	13.729 \pm 1.445	4832.392 \pm 50.014
		DeepSets	0.827 \pm 0.024	454.141 \pm 10.203	5.889 \pm 0.135	7589.795 \pm 373.293
		Transformer	0.985 \pm 0.075	274.021 \pm 1.333	6.364 \pm 0.201	4801.964 \pm 59.175

Table 17: **Variable-Dimensional.** Results for estimating the parameters of nonlinear regression models with ReLU activation function, with the expected L_2 loss according to the posterior predictive as metric.

We assume Σ_k and π_k to be known and set Σ_k to be an identity matrix and the mixing coefficients to be equal, $\pi_k = 1/K$, for all clusters k in our experiments.

D Metrics

In this section, we provide details about the metrics considered for the different tasks. We generally look at two main metrics for benchmarking performance: L_2 loss and Accuracy. For estimating the mean of a Gaussian distribution, the L_2 loss is defined as

$$GM_{L_2} = \mathbb{E}_{\mathcal{D} \sim \chi} \mathbb{E}_{\mu \sim q_\varphi(\cdot | \mathcal{D})} \left[\sum_{i=1}^{N_{\mathcal{D}}} (\mathbf{x}_i - \mu)^2 \right] \quad (44)$$

where $\mathcal{D} = \{\mathbf{x}_i\}_{i=1}^{N_{\mathcal{D}}}$. Intuitively, this captures the quality of the estimation of the mean parameter by measuring how far the observations are from it. Lower value implies better estimation of the mean parameter. Similarly, for estimating the means of a Gaussian Mixture Model, we rely on a similar metric but we also find the cluster closest to the observation, which can be defined as

$$GMM_{L_2} = \mathbb{E}_{\mathcal{D} \sim \chi} \mathbb{E}_{\mu_k \sim q_\varphi(\cdot | \mathcal{D})} \left[\sum_{i=1}^{N_{\mathcal{D}}} (\mathbf{x}_i - \mu_{\text{Match}(\mathbf{x}_i, \{\mu_1, \dots, \mu_K\})})^2 \right] \quad (45)$$

$$\text{Match}(\mathbf{x}, \{\mu_1, \dots, \mu_K\}) = \arg \min_k (\mathbf{x} - \mu_k)^2 \quad (46)$$

Objective	q_φ	Model	L_2 Loss (\downarrow)			
			1-layer		2-layers	
			1D	50D	1D	50D
Baseline	-	Random	33.972 \pm 0.273	56.891 \pm 0.213	50.843 \pm 0.270	55.383 \pm 0.124
	-	Optimization	0.330 \pm 0.000	23.257 \pm 0.029	0.672 \pm 0.007	34.900 \pm 0.120
	-	Langevin	0.296 \pm 0.003	26.292 \pm 0.266	0.356 \pm 0.006	36.598 \pm 0.354
	-	HMC	0.404 \pm 0.005	19.937 \pm 0.259	0.654 \pm 0.014	32.884 \pm 0.448
Fwd-KL	Gaussian	GRU	34.235 \pm 0.240	56.947 \pm 0.339	50.493 \pm 0.676	55.633 \pm 0.396
		DeepSets	33.979 \pm 0.273	56.900 \pm 0.213	50.844 \pm 0.284	55.390 \pm 0.132
		Transformer	33.998 \pm 0.499	56.365 \pm 0.129	49.940 \pm 0.317	55.708 \pm 0.149
Rev-KL	Gaussian	GRU	0.813 \pm 0.036	18.042 \pm 0.084	17.506 \pm 10.364	26.729 \pm 2.167
		DeepSets	0.604 \pm 0.015	22.582 \pm 2.064	24.819 \pm 0.009	28.247 \pm 0.006
		Transformer	0.896 \pm 0.097	16.726 \pm 0.724	2.249 \pm 0.651	24.442 \pm 0.643
Fwd-KL	Flow	GRU	34.989 \pm 0.475	56.847 \pm 0.289	49.915 \pm 0.980	55.505 \pm 0.479
		DeepSets	34.857 \pm 0.278	56.736 \pm 0.395	49.594 \pm 0.622	55.862 \pm 0.498
		Transformer	34.878 \pm 0.829	55.751 \pm 0.409	49.309 \pm 0.651	55.475 \pm 0.090
Rev-KL	Flow	GRU	0.969 \pm 0.038	17.454 \pm 0.041	24.796 \pm 0.028	28.258 \pm 0.008
		DeepSets	0.729 \pm 0.013	22.888 \pm 2.610	24.794 \pm 0.021	28.253 \pm 0.003
		Transformer	0.624 \pm 0.051	15.971 \pm 0.075	3.095 \pm 0.010	23.740 \pm 0.343

Table 18: **Variable-Dimensional.** Results for estimating the parameters of nonlinear regression models with TanH activation function, with the expected L_2 loss according to the posterior predictive as metric.

which intuitively captures the distance of observations from the cluster closest to them. Next, we define the metric for evaluating (non-)linear regression models as

$$(N-)LR_{L_2} = \mathbb{E}_{\mathcal{D} \sim \chi} \mathbb{E}_{\theta \sim q_\varphi(\cdot|\mathcal{D})} \left[\sum_{i=1}^{N_{\mathcal{D}}} (y_i - \text{Mode}[p(y_i|\mathbf{x}_i, \theta)])^2 \right] \quad (47)$$

Finally, for the (non-)linear classification setups, we define the accuracy metric as

$$(N-)LC_{Accuracy} = \mathbb{E}_{\mathcal{D} \sim \chi} \mathbb{E}_{\theta \sim q_\varphi(\cdot|\mathcal{D})} \left[\frac{100}{N_{\mathcal{D}}} \times \sum_{i=1}^{N_{\mathcal{D}}} \delta(y_i, \text{Mode}[p(y_i|\mathbf{x}_i, \theta)]) \right] \quad (48)$$

where $\delta(a, b) = 1$ if and only if $a = b$. Thus this metric captures the accuracy of the posterior predictive distribution. Another metric that we use to test the quality of the posterior is the symmetric KL divergence, defined as

$$\text{Symmetric } \mathbb{KL}(p(\theta|\mathcal{D}), q_\varphi(\theta|\mathcal{D})) = \frac{1}{2} \mathbb{KL}(p(\theta|\mathcal{D})||q_\varphi(\theta|\mathcal{D})) + \frac{1}{2} \mathbb{KL}(q_\varphi(\theta|\mathcal{D})||p(\theta|\mathcal{D})) \quad (49)$$

E Architecture Details

In this section, we outline the two candidate architectures that we consider for the backbone of our amortized variational inference model. We discuss the specifics of the architectures and the hyperparameters used for our experiments.

E.1 Transformer

We use a transformer model (Vaswani et al., 2017) as a permutation invariant architecture by removing positional encodings from the setup and using multiple layers of the encoder model. We append the set of

Objective	q_φ	Model	Accuracy (\uparrow)			
			1-layer		2-layers	
			2D	50D	2D	50D
Baseline	-	Random	49.951 \pm 0.287	49.904 \pm 0.281	50.040 \pm 0.467	50.044 \pm 0.239
	-	Optimization	96.762 \pm 0.034	76.139 \pm 0.033	96.810 \pm 0.009	78.225 \pm 0.129
	-	Langevin	96.077 \pm 0.027	70.142 \pm 0.229	96.564 \pm 0.113	71.328 \pm 0.410
	-	HMC	91.734 \pm 0.152	67.986 \pm 0.372	91.336 \pm 0.591	71.825 \pm 0.507
Fwd-KL	Gaussian	GRU	59.551 \pm 0.199	58.637 \pm 0.250	60.247 \pm 0.645	58.862 \pm 0.065
		DeepSets	49.946 \pm 0.285	49.910 \pm 0.282	50.032 \pm 0.466	50.040 \pm 0.242
		Transformer	59.887 \pm 0.235	58.826 \pm 0.237	60.552 \pm 0.398	58.953 \pm 0.110
Rev-KL	Gaussian	GRU	88.822 \pm 0.471	68.368 \pm 0.342	81.884 \pm 1.450	67.264 \pm 0.138
		DeepSets	91.019 \pm 0.454	61.732 \pm 0.111	82.396 \pm 0.471	67.320 \pm 0.143
		Transformer	89.988 \pm 0.197	73.744 \pm 0.319	83.399 \pm 0.841	67.167 \pm 0.028
Fwd-KL	Flow	GRU	61.179 \pm 0.833	60.225 \pm 0.115	60.400 \pm 1.019	59.027 \pm 0.212
		DeepSets	49.568 \pm 0.230	50.130 \pm 0.101	50.356 \pm 0.773	49.806 \pm 0.331
		Transformer	60.886 \pm 0.252	60.253 \pm 0.082	61.694 \pm 0.314	60.426 \pm 0.203
Rev-KL	Flow	GRU	90.363 \pm 0.709	66.197 \pm 0.118	83.443 \pm 0.619	69.053 \pm 0.256
		DeepSets	89.150 \pm 0.338	62.939 \pm 0.112	79.889 \pm 0.567	69.015 \pm 0.147
		Transformer	91.065 \pm 0.156	72.581 \pm 0.117	83.533 \pm 0.677	68.933 \pm 0.120

Table 19: **Variable-Dimensional.** Results for estimating the parameters of nonlinear classification models with ReLU activation function and two classes, with the expected accuracy according to the posterior predictive as metric.

observations with a [CLS] token before passing it to the model and use its output embedding to predict the parameters of the variational distribution. Since no positional encodings or causal masking is used in the whole setup, the final embedding of the [CLS] token becomes invariant to permutations in the set of observations, thereby leading to permutation invariance in the parameters of q_φ .

We use 4 encoder layers with a 256 dimensional attention block and 1024 feed-forward dimensions, with 4 heads in each attention block for our Transformer models to make the number of parameters comparative to the one of the DeepSets model.

E.2 DeepSets

Another framework that can process set-based input is Deep Sets (Zaheer et al., 2017). In our experiments, we used an embedding network that encodes the input into representation space, a mean aggregation operation, which ensures that the representation learned is invariant concerning the set ordering, and a regression network. The latter’s output is either used to directly parameterize a diagonal Gaussian or as conditional input to a normalizing flow, representing a summary statistics of the set input.

For DeepSets, we use 4 layers each in the embedding network and the regression network, with a mean aggregation function, ReLU activation functions, and 627 hidden dimensions to make the number of parameters comparable to those in the Transformer model.

E.3 RNN

For the recurrent neural network setup, we use the Gated Recurrent Unit (GRU). Similar to the above setups, we use a 4-layered GRU model with 256 hidden dimensions. While such an architecture is not permutation

Objective	q_φ	Model	Accuracy (\uparrow)			
			1-layer		2-layers	
			2D	50D	2D	50D
Baseline	-	Random	19.847 \pm 0.352	20.052 \pm 0.066	19.874 \pm 0.222	20.028 \pm 0.106
	-	Optimization	94.607 \pm 0.011	56.091 \pm 0.158	93.873 \pm 0.028	60.253 \pm 0.053
	-	Langevin	90.815 \pm 0.341	46.072 \pm 0.225	91.849 \pm 0.088	49.808 \pm 0.337
	-	HMC	81.145 \pm 0.303	44.561 \pm 0.309	79.559 \pm 0.483	50.967 \pm 0.436
Fwd-KL	Gaussian	GRU	30.960 \pm 0.638	32.164 \pm 0.151	31.017 \pm 0.397	32.224 \pm 0.108
		DeepSets	19.846 \pm 0.348	20.053 \pm 0.063	19.871 \pm 0.226	20.032 \pm 0.104
		Transformer	30.652 \pm 0.401	32.208 \pm 0.178	31.148 \pm 0.429	32.342 \pm 0.217
Rev-KL	Gaussian	GRU	72.874 \pm 0.113	37.987 \pm 0.119	56.999 \pm 0.599	29.971 \pm 0.248
		DeepSets	69.456 \pm 0.370	36.712 \pm 0.249	55.193 \pm 0.538	36.417 \pm 9.779
		Transformer	73.531 \pm 0.391	44.702 \pm 0.165	57.724 \pm 0.332	30.175 \pm 0.177
Fwd-KL	Flow	GRU	33.232 \pm 0.607	33.704 \pm 0.026	31.937 \pm 0.483	32.370 \pm 0.284
		DeepSets	19.898 \pm 0.156	19.950 \pm 0.256	20.062 \pm 0.347	20.064 \pm 0.207
		Transformer	32.916 \pm 0.194	33.766 \pm 0.137	32.374 \pm 0.301	33.846 \pm 0.486
Rev-KL	Flow	GRU	77.997 \pm 0.663	38.715 \pm 0.153	61.947 \pm 0.294	51.962 \pm 0.812
		DeepSets	68.957 \pm 0.551	37.123 \pm 0.108	51.145 \pm 14.201	42.707 \pm 12.882
		Transformer	77.867 \pm 2.241	44.156 \pm 0.485	57.410 \pm 0.088	52.077 \pm 0.077

Table 20: **Fixed-Dimensional**. Results for estimating the parameters of nonlinear classification models with ReLU activation function and five classes, with the expected accuracy according to the posterior predictive as metric.

invariant, by training on tasks that require such invariance could encourage learning of solution structure that respects this invariance.

E.4 Normalizing Flows

Assuming a Gaussian posterior distribution as the approximate often leads to poor results as the true posterior distribution can be far from the Gaussian shape. To allow for more flexible posterior distributions, we use normalizing flows (Kingma & Dhariwal, 2018; Kobyzev et al., 2020; Papamakarios et al., 2021; Rezende & Mohamed, 2015) for approximating $q_\varphi(\boldsymbol{\theta}|\mathcal{D})$ conditioned on the output of the summary network h_ψ . Specifically, let $g_\nu : \mathbf{z} \mapsto \boldsymbol{\theta}$ be a diffeomorphism parameterized by a conditional invertible neural network (cINN) with network parameters ν such that $\boldsymbol{\theta} = g_\nu(\mathbf{z}; h_\psi(\mathcal{D}))$. With the change-of-variables formula it follows that $p(\boldsymbol{\theta}) = p(\mathbf{z}) |\det \frac{\partial}{\partial \mathbf{z}} g_\nu(\mathbf{z}; h_\psi(\mathcal{D}))|^{-1} = p(\mathbf{z}) |\det J_\nu(\mathbf{z}; h_\psi(\mathcal{D}))|^{-1}$, where J_ν is the Jacobian matrix of g_ν . Further, integration by substitution gives us $d\boldsymbol{\theta} = |\det J_\nu(\mathbf{z}; h_\psi(\mathcal{D}))| d\mathbf{z}$ to rewrite the objective from eq. 12 as:

$$\arg \min_{\varphi} \mathbb{KL}[q_\varphi(\boldsymbol{\theta}|\mathcal{D})||p(\boldsymbol{\theta}|\mathcal{D})] \quad (50)$$

$$= \arg \min_{\varphi} \mathbb{E}_{\mathcal{D} \sim \chi} \mathbb{E}_{\boldsymbol{\theta} \sim q_\varphi(\boldsymbol{\theta}|\mathcal{D})} [\log q_\varphi(\boldsymbol{\theta}|\mathcal{D}) - \log p(\boldsymbol{\theta}, \mathcal{D})] \quad (51)$$

$$= \arg \min_{\varphi=\{\psi, \nu\}} \mathbb{E}_{\mathcal{D} \sim \chi} \mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})} \left[\log \frac{q_\nu(\mathbf{z}|h_\psi(\mathcal{D}))}{|\det J_\nu(\mathbf{z}; h_\psi(\mathcal{D}))|} - \log p(g_\nu(\mathbf{z}; h_\psi(\mathcal{D})), \mathcal{D}) \right] \quad (52)$$

As shown in BayesFlow (Radev et al., 2020), the normalizing flow g_ν and the summary network h_ψ can be trained simultaneously. The AllInOneBlock coupling block architecture of the FrEIA Python package (Ardizzone et al., 2018), which is very similar to the RNVP style coupling block (Dinh et al., 2017), is used

Objective	q_φ	Model	Accuracy (\uparrow)			
			1-layer		2-layers	
			2D	50D	2D	50D
Baseline	-	Random	50.278 \pm 0.337	50.028 \pm 0.064	50.188 \pm 0.479	49.982 \pm 0.084
	-	Optimization	96.943 \pm 0.012	68.086 \pm 0.016	94.444 \pm 0.024	63.950 \pm 0.022
	-	Langevin	95.143 \pm 0.094	61.694 \pm 0.340	92.719 \pm 0.016	57.447 \pm 0.437
	-	HMC	92.489 \pm 0.338	59.963 \pm 0.202	87.548 \pm 0.094	56.319 \pm 0.882
Fwd-KL	Gaussian	GRU	50.274 \pm 0.337	50.023 \pm 0.060	50.187 \pm 0.477	49.992 \pm 0.072
		DeepSets	50.271 \pm 0.334	50.024 \pm 0.061	50.188 \pm 0.472	49.984 \pm 0.077
		Transformer	50.273 \pm 0.336	50.031 \pm 0.066	50.191 \pm 0.471	49.994 \pm 0.078
Rev-KL	Gaussian	GRU	89.270 \pm 0.272	50.014 \pm 0.048	50.191 \pm 0.463	49.995 \pm 0.080
		DeepSets	89.788 \pm 0.213	50.018 \pm 0.061	50.191 \pm 0.478	49.977 \pm 0.075
		Transformer	89.366 \pm 0.108	64.926 \pm 0.260	50.182 \pm 0.469	49.986 \pm 0.071
Fwd-KL	Flow	GRU	49.651 \pm 0.040	50.119 \pm 0.068	49.987 \pm 0.015	49.904 \pm 0.018
		DeepSets	49.639 \pm 0.031	50.113 \pm 0.065	49.988 \pm 0.022	49.910 \pm 0.043
		Transformer	49.636 \pm 0.040	50.115 \pm 0.063	49.989 \pm 0.017	49.909 \pm 0.042
Rev-KL	Flow	GRU	49.769 \pm 0.141	50.082 \pm 0.091	49.915 \pm 0.070	50.004 \pm 0.084
		DeepSets	49.782 \pm 0.073	50.080 \pm 0.087	49.831 \pm 0.152	49.994 \pm 0.078
		Transformer	63.233 \pm 19.243	50.026 \pm 0.047	49.869 \pm 0.207	50.036 \pm 0.056

Table 21: **Variable-Dimensional.** Results for estimating the parameters of nonlinear classification models with TanH activation function and two classes, with the expected accuracy according to the posterior predictive as metric.

as the basis for the cINN. AllInOneBlock combines the most common architectural components, such as ActNorm, permutation, and affine coupling operations.

For our experiments, 6 coupling blocks define the normalizing flow network, each with a 1 hidden-layered non-linear feed-forward subnetwork with ReLU non-linearity and 128 hidden dimensions.

F Experimental Details

Unless specified, we obtain a stream of datasets for all our experiments by simply sampling from the assumed probabilistic model, where the number of observations n is sampled uniformly in the range [64, 128]. For efficient mini-batching over datasets with different cardinalities, we sample datasets with maximum cardinality (128) and implement different cardinalities by masking out different numbers of observations for different datasets whenever required.

To evaluate both our proposed approach and the baselines, we compute an average of the predictive performances across 25 different posterior samples for each of the 100 fixed test datasets for all our experiments. That means for our proposed approach, we sample 25 different parameter vectors from the approximate posterior that we obtain. For MCMC, we rely on 25 MCMC samples, and for optimization, we train 25 different parameter vectors where the randomness comes from initialization. For the optimization baseline, we perform a quick hyperparameter search over the space $\{0.01, 0.003, 0.001, 0.0003, 0.0001, 0.00003\}$ to pick the best learning rate that works for all of the test datasets and then use it to train for 1000 iterations using the Adam optimizer (Kingma & Ba, 2014). For the MCMC baseline, we use the open-sourced implementation of Langevin-based MCMC sampling² where we leave a chunk of the starting samples as burn-in and then

²<https://github.com/alisiahkoohi/Langevin-dynamics>

Objective	q_φ	Model	Accuracy (\uparrow)			
			1-layer		2-layers	
			2D	50D	2D	50D
Baseline	-	Random	20.041 \pm 0.136	20.002 \pm 0.079	19.914 \pm 0.111	19.954 \pm 0.005
	-	Optimization	92.059 \pm 0.012	40.722 \pm 0.027	88.848 \pm 0.005	34.136 \pm 0.041
	-	Langevin	88.357 \pm 0.309	30.941 \pm 0.097	83.788 \pm 0.170	26.538 \pm 0.206
	-	HMC	79.161 \pm 0.292	27.508 \pm 0.379	74.987 \pm 0.130	25.377 \pm 0.246
Fwd-KL	Gaussian	GRU	20.264 \pm 0.139	20.158 \pm 0.056	20.215 \pm 0.138	20.093 \pm 0.028
		DeepSets	20.042 \pm 0.133	20.000 \pm 0.088	19.916 \pm 0.114	19.955 \pm 0.007
		Transformer	20.240 \pm 0.126	20.153 \pm 0.070	20.118 \pm 0.131	20.089 \pm 0.022
Rev-KL	Gaussian	GRU	66.565 \pm 7.725	20.011 \pm 0.099	19.913 \pm 0.125	19.954 \pm 0.022
		DeepSets	57.294 \pm 0.362	20.011 \pm 0.091	19.915 \pm 0.115	19.959 \pm 0.020
		Transformer	72.865 \pm 1.340	22.185 \pm 1.872	19.911 \pm 0.124	19.953 \pm 0.014
Fwd-KL	Flow	GRU	19.963 \pm 0.239	20.176 \pm 0.056	19.952 \pm 0.189	20.156 \pm 0.079
		DeepSets	19.757 \pm 0.259	20.045 \pm 0.071	19.692 \pm 0.184	20.019 \pm 0.069
		Transformer	19.925 \pm 0.262	20.185 \pm 0.059	19.882 \pm 0.175	20.159 \pm 0.070
Rev-KL	Flow	GRU	67.042 \pm 2.230	20.065 \pm 0.060	19.707 \pm 0.245	19.989 \pm 0.101
		DeepSets	35.220 \pm 10.870	20.000 \pm 0.054	19.739 \pm 0.216	19.966 \pm 0.019
		Transformer	72.798 \pm 1.049	20.032 \pm 0.040	19.752 \pm 0.276	20.017 \pm 0.037

Table 22: **Variable-Dimensional.** Results for estimating the parameters of nonlinear classification models with TanH activation function and five classes, with the expected accuracy according to the posterior predictive as metric.

start accepting samples after a regular interval (to not make them correlated). The details about the burn-in time and the regular interval for acceptance are provided in the corresponding experiments’ sections below.

For our proposed approach of amortized inference, we do not consider explicit hyperparameter optimization and simply use a learning rate of 1e-4 with the Adam optimizer. For all experiments, we used linear scaling of the KL term in the training objectives as described in (Higgins et al., 2017), which we refer to as warmup. Furthermore, training details for each experiment can be found below.

F.1 Fixed-Dim

In this section, we provide the experimental details relevant to reproducing the results of Section 4. All the models are trained with streaming data from the underlying probabilistic model, such that every iteration of training sees a new set of datasets. Training is done with a batch size of 128, representing the number of datasets seen during one optimization step. Evaluations are done with 25 samples and we ensure that the test datasets used for each probabilistic model are the same across all the compared methods, i.e., baselines, forward KL, and reverse KL. We train the amortized inference model and the forward KL baselines for the following different probabilistic models:

Mean of Gaussian (GM): We train the amortization models over 20,000 iterations for both the 2-dimensional as well as the 100-dimensional setup. We use a linear warmup with 5000 iterations over which the weight of the KL term in our proposed approach scales linearly from 0 to 1. We use an identity covariance matrix for the data-generating process, but it can be easily extended to the case of correlated or diagonal covariance-based Gaussian distributions.

Gaussian Mixture Model (GMM): We train the mixture model setup for 200,000 iterations with 50,000 iterations of warmup. We mainly experiment with 2-dimensional and 5-dimensional mixture models,

Objective	q_φ	Model	L_2 Loss (\downarrow)						$\leftarrow \chi_{real}$	
			Linear Model MLP-TanH Data			MLP-TanH Model Linear Data				$\leftarrow \chi_{sim}$
			LR	NLR	GP	LR	NLR	GP		
Baseline	-	Random	-	17.761 \pm 0.074	-	17.847 \pm 0.355	-	-		
	-	Optimization	-	1.213 \pm 0.000	-	0.360 \pm 0.001	-	-		
	-	Langevin	-	1.218 \pm 0.002	-	0.288 \pm 0.001	-	-		
	-	HMC	-	1.216 \pm 0.002	-	0.275 \pm 0.001	-	-		
Fwd-KL	Gaussian	GRU	2.415 \pm 0.269	-	-	-	15.632 \pm 0.283	-		
		DeepSets	1.402 \pm 0.017	-	-	-	16.046 \pm 0.393	-		
		Transformer	2.216 \pm 0.097	-	-	-	15.454 \pm 0.246	-		
Rev-KL	Gaussian	GRU	1.766 \pm 0.044	1.216 \pm 0.001	4.566 \pm 0.199	0.375 \pm 0.001	0.386 \pm 0.002	0.524 \pm 0.019		
		DeepSets	1.237 \pm 0.006	1.216 \pm 0.001	3963.694 \pm 5602.411	0.365 \pm 0.000	0.377 \pm 0.003	0.385 \pm 0.011		
		Transformer	1.892 \pm 0.113	1.226 \pm 0.001	4.313 \pm 0.707	0.367 \pm 0.006	0.382 \pm 0.003	0.458 \pm 0.048		
Fwd-KL	Flow	GRU	2.180 \pm 0.024	-	-	-	9.800 \pm 0.473	-		
		DeepSets	1.713 \pm 0.244	-	-	-	15.253 \pm 0.403	-		
		Transformer	1.632 \pm 0.070	-	-	-	7.949 \pm 0.419	-		
Rev-KL	Flow	GRU	1.830 \pm 0.081	1.214 \pm 0.001	5.690 \pm 0.196	0.346 \pm 0.004	0.349 \pm 0.001	0.520 \pm 0.015		
		DeepSets	1.282 \pm 0.036	1.218 \pm 0.001	11.690 \pm 10.602	0.339 \pm 0.003	0.344 \pm 0.002	0.397 \pm 0.026		
		Transformer	1.471 \pm 0.016	1.226 \pm 0.004	5.194 \pm 0.320	0.346 \pm 0.002	0.347 \pm 0.001	0.480 \pm 0.030		

Table 23: **Model Misspecification.** Results for model misspecification under different training data χ_{sim} , when evaluated under MLP-TanH and Linear Data (χ_{real}), with the underlying model as a linear and MLP-TanH model respectively.

with 2 and 5 mixture components for each setup. While we do use an identity covariance matrix for the data-generating process, again, it can be easily extended to other cases.

Linear Regression (LR): The amortization models for this setup are trained for 50,000 iterations with 12,500 iterations of warmup. The feature dimensions considered for this task are 1 and 100 dimensions, and the predictive variance σ^2 is assumed to be known and set as 0.25.

Nonlinear Regression (NLR): We train the setup for 100,000 iterations with 25,000 iterations consisting of warmup. The feature dimensionalities considered are 1-dimensional and 25-dimensional, and training is done with a known predictive variance similar to the LR setup. For the probabilistic model, we consider both a 1-layered and a 2-layered multi-layer perceptron (MLP) network with 32 hidden units in each, and either a RELU or TANH activation function.

Linear Classification (LC): We experiment with 2-dimensional and 100-dimensional setups with training done for 50,000 iterations, out of which 12,500 are used for warmup. Further, we train for both binary classification as well as a 5-class classification setup.

Nonlinear Classification (NLC): We experiment with 2-dimensional and 25-dimensional setups with training done for 100,000 iterations, out of which 2,5000 are used for warmup. Further, we train for both binary classification as well as a 5-class classification setup. For the probabilistic model, we consider both a 1-layered and a 2-layered multi-layer perceptron (MLP) network with 32 hidden units in each, and either a RELU or TANH activation function.

F.2 Variable-Dim

In this section, we provide the experimental details relevant to reproducing the results of Section 4. All the models are trained with streaming data from the underlying probabilistic model, such that every iteration of training sees a new set of datasets. Training is done with a batch size of 128, representing the number of datasets seen during one optimization step. Further, we ensure that the datasets sampled resemble a

Objective	q_φ	Model	L_2 Loss (\downarrow)						$\leftarrow \chi_{real}$	
			Linear Model GP Data			MLP-TanH Model GP Data				$\leftarrow \chi_{sim}$
			LR	NLR	GP	LR	NLR	GP		
Baseline	-	Random	-	-	2.681 \pm 0.089	-	-	16.236 \pm 0.381		
	-	Optimization	-	-	0.263 \pm 0.000	-	-	0.007 \pm 0.000		
	-	Langevin	-	-	0.266 \pm 0.001	-	-	0.022 \pm 0.001		
	-	HMC	-	-	0.266 \pm 0.000	-	-	0.090 \pm 0.002		
Fwd-KL	Gaussian	GRU	0.268 \pm 0.000	-	-	-	14.077 \pm 0.368	-		
		DeepSets	0.269 \pm 0.001	-	-	-	14.756 \pm 0.280	-		
		Transformer	0.270 \pm 0.001	-	-	-	14.733 \pm 0.513	-		
Rev-KL	Gaussian	GRU	0.268 \pm 0.000	0.269 \pm 0.000	0.266 \pm 0.000	0.334 \pm 0.005	0.157 \pm 0.003	0.080 \pm 0.003		
		DeepSets	0.269 \pm 0.000	0.269 \pm 0.000	0.265 \pm 0.000	0.331 \pm 0.003	0.146 \pm 0.002	0.063 \pm 0.000		
		Transformer	0.269 \pm 0.000	0.269 \pm 0.000	0.267 \pm 0.000	0.310 \pm 0.013	0.155 \pm 0.006	0.066 \pm 0.004		
Fwd-KL	Flow	GRU	0.268 \pm 0.000	-	-	-	9.756 \pm 0.192	-		
		DeepSets	0.269 \pm 0.001	-	-	-	14.345 \pm 0.628	-		
		Transformer	0.269 \pm 0.000	-	-	-	8.557 \pm 0.561	-		
Rev-KL	Flow	GRU	0.268 \pm 0.000	0.270 \pm 0.001	0.266 \pm 0.000	0.289 \pm 0.011	0.120 \pm 0.004	0.059 \pm 0.003		
		DeepSets	0.269 \pm 0.000	0.269 \pm 0.001	0.266 \pm 0.000	0.270 \pm 0.008	0.115 \pm 0.002	0.059 \pm 0.002		
		Transformer	0.269 \pm 0.001	0.270 \pm 0.000	0.267 \pm 0.000	0.293 \pm 0.008	0.120 \pm 0.005	0.055 \pm 0.002		

Table 24: **Model Misspecification.** Results for model misspecification under different training data χ_{sim} , when evaluated under GP Data (χ_{real}), with the underlying model as a linear and MLP-TanH model respectively.

uniform distribution over the feature dimensions, ranging from 1-dimensional to the maximal dimensional setup. Evaluations are done with 25 samples and we ensure that the test datasets used for each probabilistic model are the same across all the compared methods, i.e., baselines, forward KL, and reverse KL. We train the amortized inference model and the forward KL baselines for the following different probabilistic models:

Mean of Gaussian (GM): We train the amortization models over 50,000 iterations using a linear warmup with 12,500 iterations over which the weight of the KL term in our proposed approach scales linearly from 0 to 1. We use an identity covariance matrix for the data-generating process, but it can be easily extended to the case of correlated or diagonal covariance-based Gaussian distributions. In this setup, we consider a maximum of 100 feature dimensions.

Gaussian Mixture Model (GMM): We train the mixture model setup for 500,000 iterations with 125,000 iterations of warmup. We set the maximal feature dimensions as 5 and experiment with 2 and 5 mixture components. While we do use an identity covariance matrix for the data-generating process, again, it can be easily extended to other cases.

Linear Regression (LR): The amortization models for this setup are trained for 100,000 iterations with 25,000 iterations of warmup. The maximal feature dimension considered for this task is 100-dimensional, and the predictive variance σ^2 is assumed to be known and set as 0.25.

Nonlinear Regression (NLR): We train the setup for 250,000 iterations with 62,500 iterations consisting of warmup. The maximal feature dimension considered is 100-dimensional, and training is done with a known predictive variance similar to the LR setup. For the probabilistic model, we consider both a 1-layered and a 2-layered multi-layer perceptron (MLP) network with 32 hidden units in each, and either a RELU or TANH activation function.

Linear Classification (LC): We experiment with a maximal 100-dimensional setup with training done for 100,000 iterations, out of which 25,000 are used for warmup. Further, we train for both binary classification as well as a 5-class classification setup.

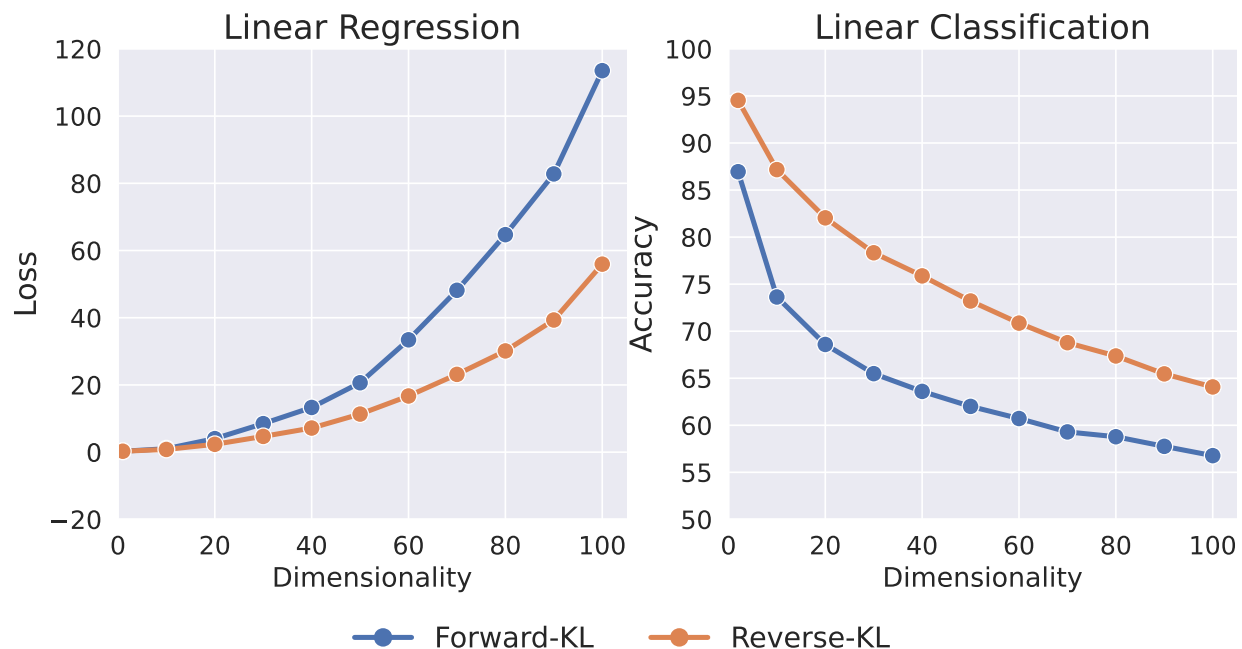


Figure 5: **Trends of Performance over different Dimensions in Variable Dimensionality Setup:** We see that our proposed reverse KL methodology outperforms the forward KL one.

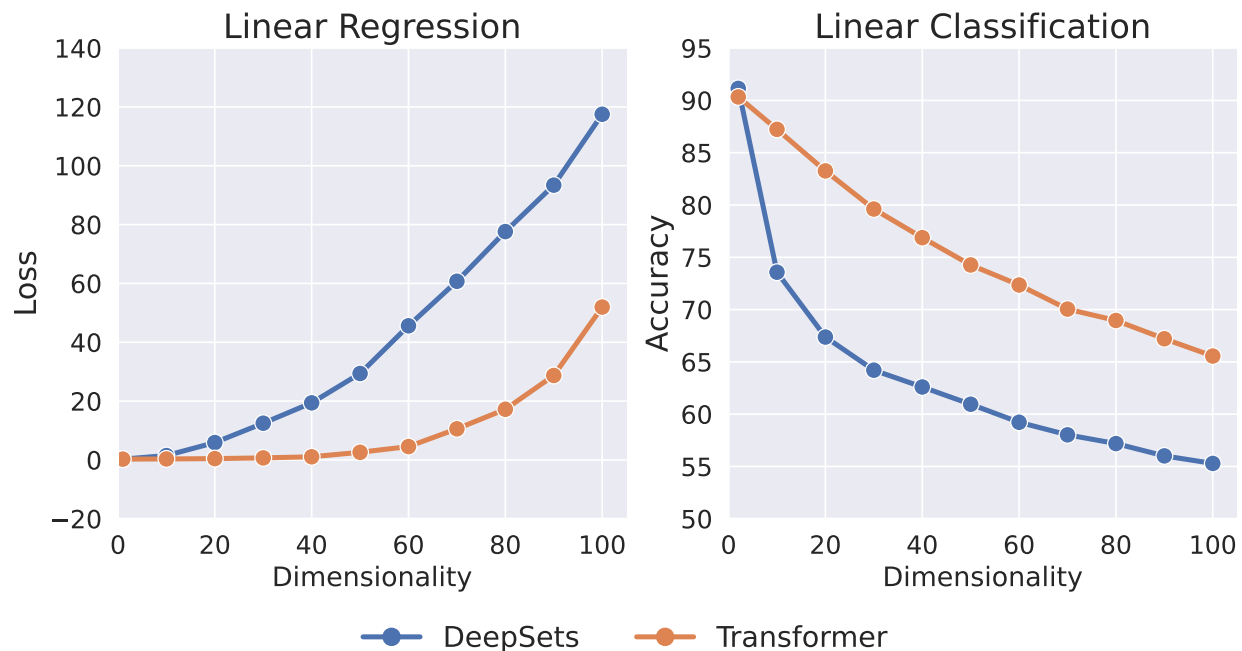


Figure 6: **Trends of Performance over different Dimensions in Variable Dimensionality Setup:** We see that transformer models generalize better to different dimensional inputs than DeepSets.

Nonlinear Classification (NLC): We experiment with a maximal 100-dimensional setup with training done for 250,000 iterations, out of which 62,500 are used for warmup. Further, we train for both binary classification as well as a 5-class classification setup. For the probabilistic model, we consider both a 1-layered and a 2-layered multi-layer perceptron (MLP) network with 32 hidden units in each, and either a RELU or TANH activation function.

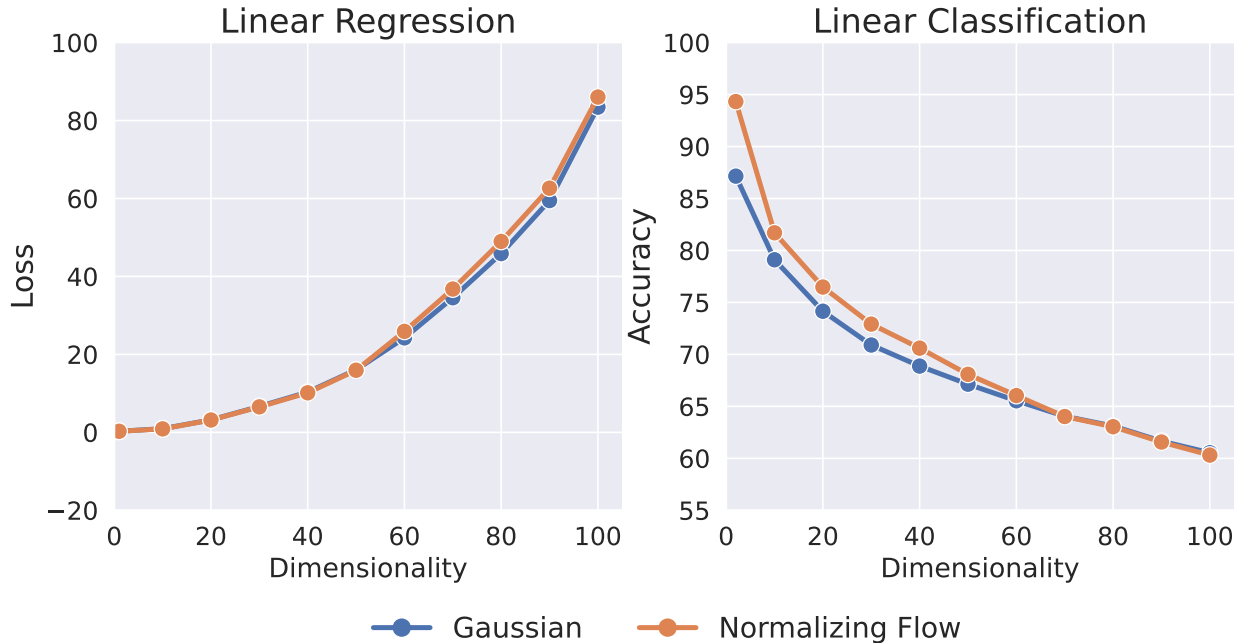


Figure 7: **Trends of Performance over different Dimensions in Variable Dimensionality Setup:** We see that normalizing flows leads to similar performances than Gaussian based variational approximation.

F.3 Model Misspecification

In this section, we provide the experimental details relevant to reproducing the results of Section 4. All models during this experiment are trained with streaming data from the currently used dataset-generating function χ , such that every iteration of training sees a new batch of datasets. Training is done with a batch size of 128, representing the number of datasets seen during one optimization step. Evaluation for all models is done with 10 samples from each dataset-generator used in the respective experimental subsection and we ensure that the test datasets are the same across all compared methods, i.e., baselines, forward KL, and reverse KL.

Linear Regression Model: The linear regression amortization models are trained following the training setting for linear regression fixed dimensionality, that is, 50,000 training iterations with 12,500 iterations of warmup. The feature dimension considered for this task is 1-dimension. The model is trained separately on datasets from three different generators χ : linear regression, nonlinear regression, and Gaussian processes, and evaluated after training on test datasets from all of them. For training with datasets from the linear regression probabilistic model, the predictive variance σ^2 is assumed to be known and set as 0.25. The same variance is used for generating datasets from the nonlinear regression dataset generator with 1 layer, 32 hidden units, and TANH activation function. Lastly, datasets from the Gaussian process-based generator are sampled similarly, using the GPytorch library Gardner et al. (2018), where datasets are sampled of varying cardinality, ranging from 64 to 128. We use a zero-mean Gaussian Process (GP) with a unit lengthscale radial-basis function (RBF) kernel serving as the covariance matrix. Further, we use a very small noise of $\sigma^2 = 1e^{-6}$ in the likelihood term of the GP. Forward KL training in this experiment can only be done when the amortization model and the dataset-generating function are the same: when we train on datasets from the linear regression-based χ . Table 23 provides a detailed overview of the results.

Nonlinear Regression Models: The nonlinear regression amortization models are trained following the training setting for nonlinear regression fixed dimensionality, that is, 100,000 training iterations with 25,000 iterations of warmup. Here, we consider two single-layer perceptions with 32 hidden units with a TANH activation function. The feature dimensionality considered is 1 dimension. We consider the same dataset-generating functions as in the misspecification experiment for a linear regression model above. However, the activation function used in the nonlinear regression dataset generator matches the activation function of

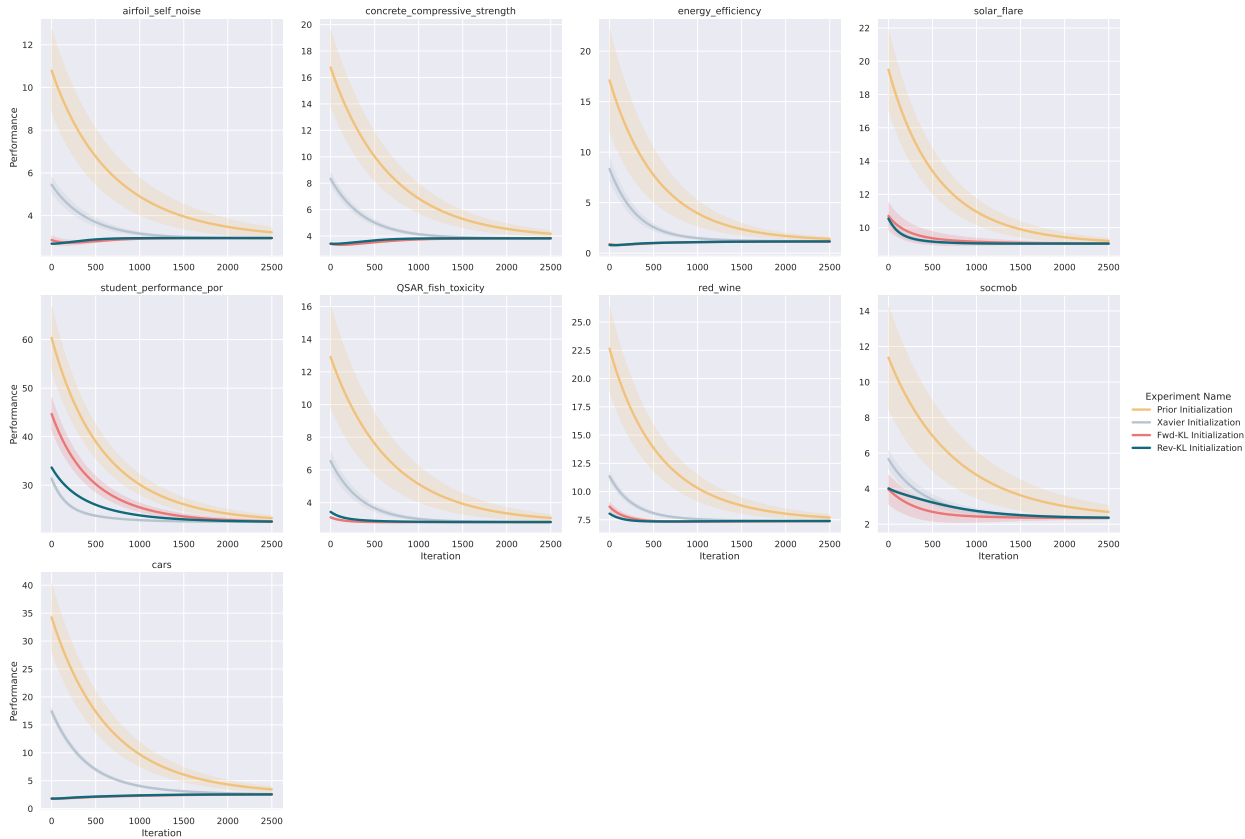


Figure 8: **Tabular Experiments | Linear Regression with Diagonal Gaussian:** For every regression dataset from the OpenML platform considered, we initialize the parameters of a linear regression-based probabilistic model with the amortized inference models which were trained with a diagonal Gaussian assumption. The parameters are then further trained with maximum-a-posteriori (MAP) estimate with gradient descent. Reverse and Forward KL denote initialization with the correspondingly trained amortized model. Prior refers to a MAP-based optimization baseline initialized from the prior $\mathcal{N}(\mathbf{0}, I)$, whereas Xavier refers to initialization from the Xavier initialization scheme.

the currently trained amortization model. In this case, forward KL training is possible in the two instances when trained on datasets from the corresponding nonlinear regression probabilistic model. A more detailed overview of the results can be found in Table 23 and 24.

F.4 Tabular Experiments

For the tabular experiments, we train the amortized inference models for (non-)linear regression (NLR/LR) as well as (non-)linear classification (NLC/LC) with $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ as opposed to $\mathbf{x} \sim \mathcal{U}(-\mathbf{1}, \mathbf{1})$ in the dataset generating process χ , with the rest of the settings the same as MAXIMUM-DIM experiments. For the nonlinear setups, we only consider the RELU case as it has seen predominant success in deep learning. Further, we only consider a 1-hidden layer neural network with 32 hidden dimensions in the probabilistic model.

After having trained the amortized inference models, both for forward and reverse KL setups, we evaluate them on real-world tabular datasets. We first collect a subset of tabular datasets from the OpenML platform as outlined in Appendix G. Then, for each dataset, we perform a 5-fold cross-validation evaluation where the dataset is chunked into 5 bins, of which, at any time, 4 are used for training and one for evaluation. This procedure is repeated five times so that every chunk is used for evaluation once.

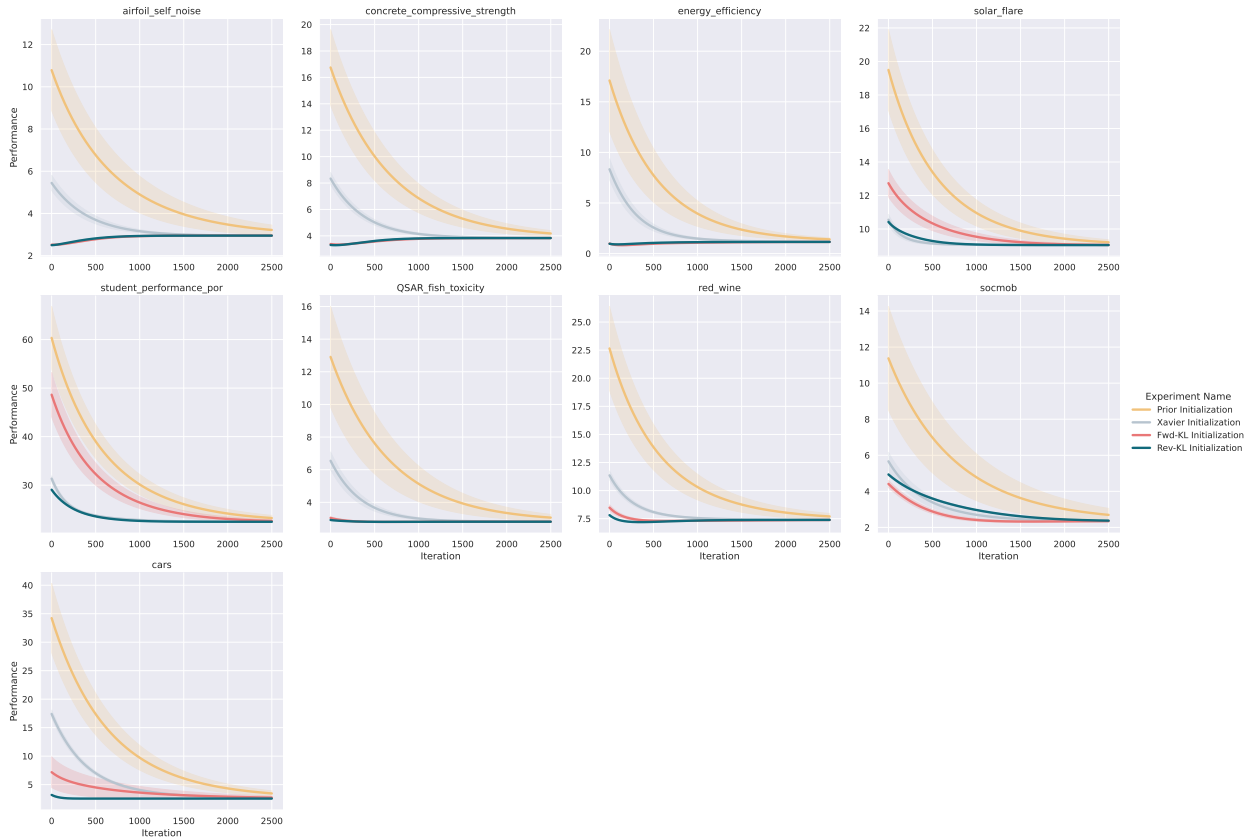


Figure 9: **Tabular Experiments | Linear Regression with Normalizing Flow:** For every regression dataset from the OpenML platform considered, we initialize the parameters of a linear regression-based probabilistic model with the amortized inference models which were trained with a normalizing flow-based model. The parameters are then further trained with maximum-a-posteriori (MAP) estimate with gradient descent. Reverse and Forward KL denote initialization with the correspondingly trained amortized model. Prior refers to a MAP-based optimization baseline initialized from the prior $\mathcal{N}(0, I)$, whereas Xavier refers to initialization from the Xavier initialization scheme.

For each dataset, we normalize the observations and the targets so that they have zero mean and unit standard deviation. For the classification setups, we only normalize the inputs as the targets are categorical. For both forward KL and reverse KL amortization models, we initialize the probabilistic model from samples from the amortized model and then further finetune it via dataset-specific maximum a posteriori optimization. We repeat this setup over 25 different samples from the inference model. In contrast, for the optimization baseline, we initialize the probabilistic models’ parameters from $\mathcal{N}(0, I)$, which is the prior that we consider, and then train 25 such models with maximum a posteriori objective using Adam optimizer.

While we see that the amortization models, particularly the reverse KL model, lead to much better initialization and convergence, it is important to note that the benefits vanish if we initialize using the Xavier-init initialization scheme. However, we believe that this is not a fair comparison as it means that we are considering a different prior now, while the amortized models were trained with $\mathcal{N}(0, I)$ prior. We defer the readers to the section below for additional discussion and experimental results.

G OpenML Datasets

For the tabular regression problems, we consider the suite of tasks outlined in *OpenML-CTR23 - A curated tabular regression benchmarking suite* (Fischer et al., 2023), from which we further filter out datasets that

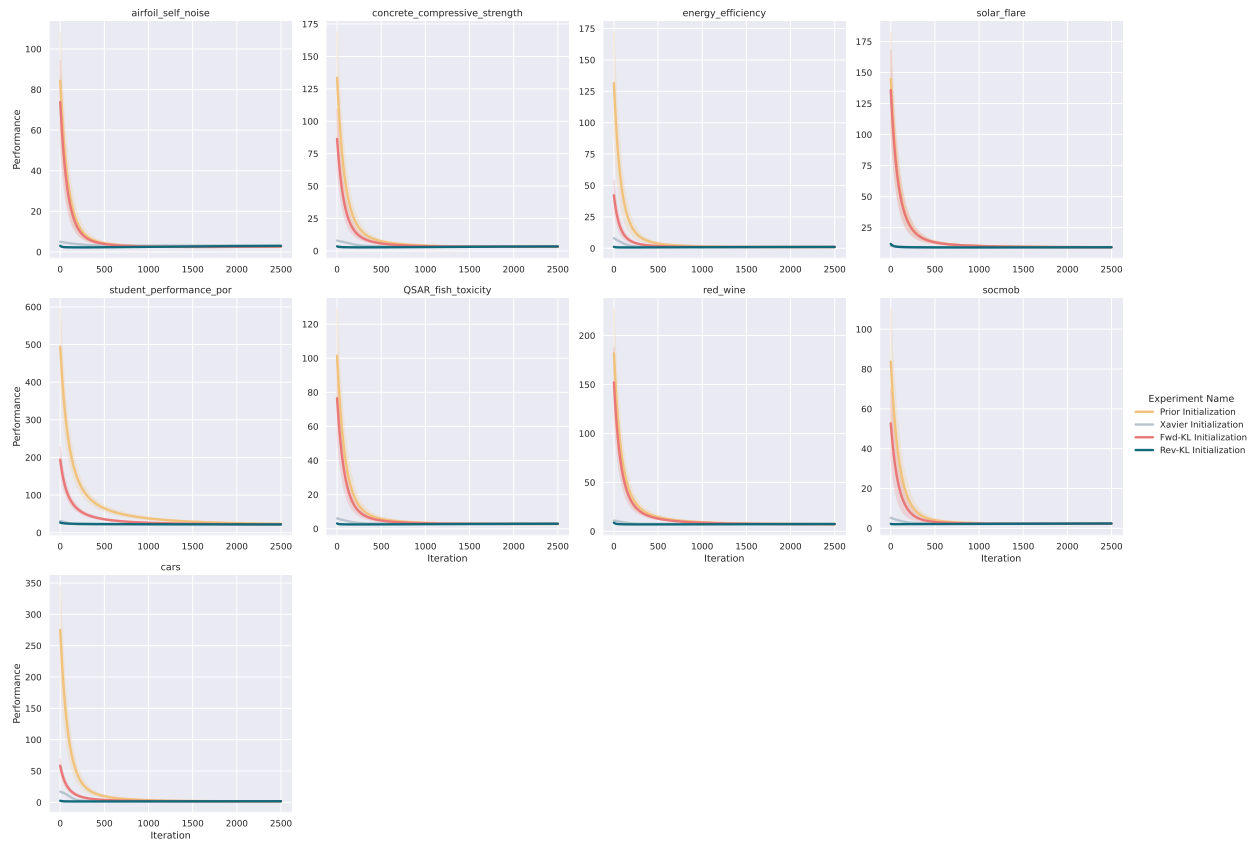


Figure 10: **Tabular Experiments | Nonlinear Regression with Diagonal Gaussian:** For every regression dataset from the OpenML platform considered, we initialize the parameters of a nonlinear regression-based probabilistic model with the amortized inference models which were trained with a diagonal Gaussian assumption. The parameters are then further trained with maximum-a-posteriori (MAP) estimate with gradient descent. Reverse and Forward KL denote initialization with the correspondingly trained amortized model. Prior refers to a MAP-based optimization baseline initialized from the prior $\mathcal{N}(0, I)$, whereas Xavier refers to initialization from the Xavier initialization scheme.

have more than 2000 examples and 100 features. We also remove datasets with missing information and NaNs. Similarly, we consider the *OpenML-CC18 Curated Classification benchmark* (Bischl et al., 2019) suite of tasks for classification and perform a similar filtering algorithm. We remove datasets with missing information and NaNs, as well as datasets with more than 2000 examples and 100 features. In addition, we also exclude datasets that are not made for binary classification. At the end of this filtering mechanism, we end up with 9 regression and 13 classification problems, and our dataset filtration pipeline is heavily inspired by Hollmann et al. (2022). We provide the datasets considered for both regression and classification below:

Regression: AIRFOIL_SELF_NOISE, CONCRETE_COMPRESSIVE_STRENGTH, ENERGY_EFFICIENCY, SOLAR_FLARE, STUDENT_PERFORMANCE_POR, QSAR_FISH_TOXICITY, RED_WINE, SOCMOB and CARS.

Classification: CREDIT-G, DIABETES, TIC-TAC-TOE, PC4, PC3, KC2, PC1, BANKNOTE-AUTHENTICATION, BLOOD-TRANSFUSION-SERVICE-CENTER, ILPD, QSAR-BIODEG, WDBC and CLIMATE-MODEL-SIMULATION-CRASHES.

H Additional Experiments

In this section, we outline the additional experiments we conducted in obtaining Bayesian posteriors for the different probabilistic models for different hyperparameters and their downstream uses. We provide a comprehensive account of the results in the relevant sections below.

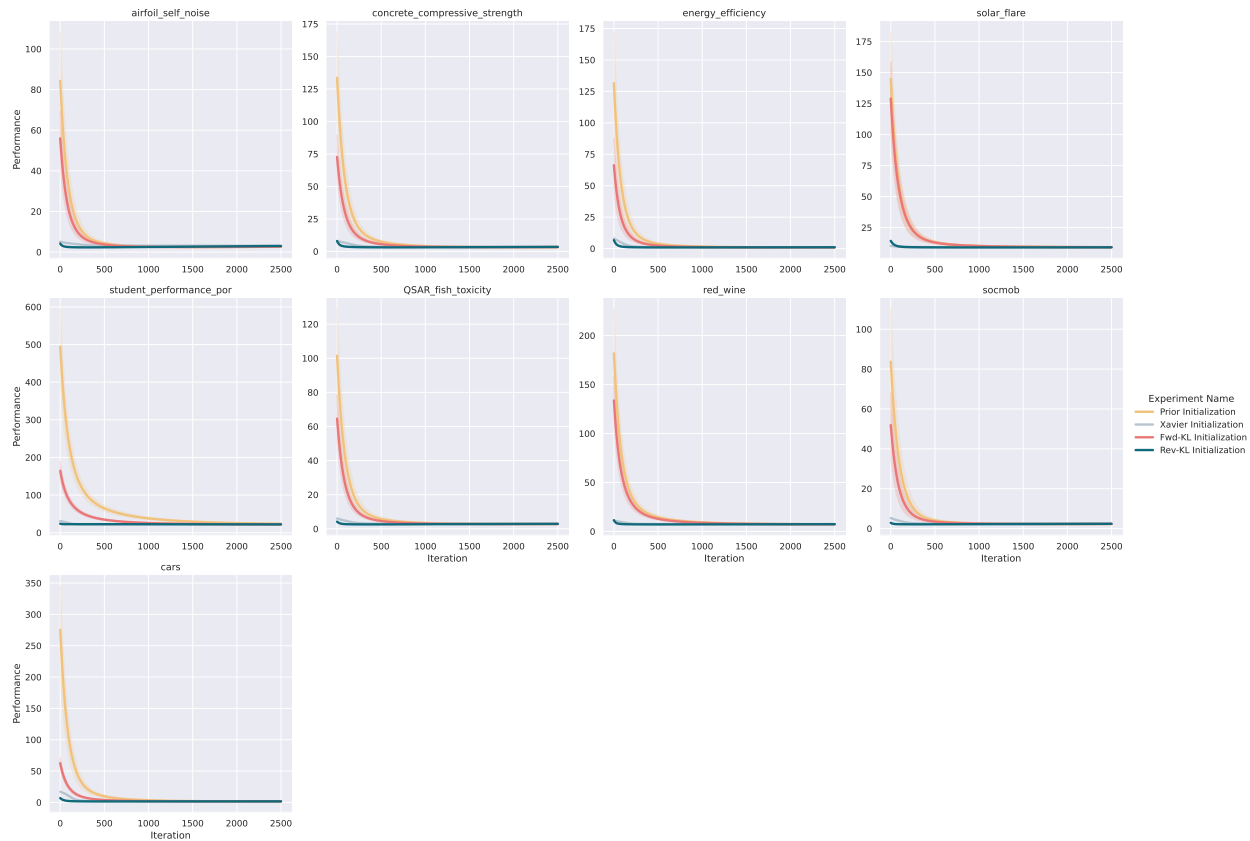


Figure 11: **Tabular Experiments | Nonlinear Regression with Normalizing Flow**: For every regression dataset from the OpenML platform considered, we initialize the parameters of a nonlinear regression-based probabilistic model with the amortized inference models which were trained with a normalizing flow-based model. The parameters are then further trained with maximum-a-posteriori (MAP) estimate with gradient descent. Reverse and Forward KL denote initialization with the correspondingly trained amortized model. Prior refers to a MAP-based optimization baseline initialized from the prior $\mathcal{N}(0, I)$, whereas Xavier refers to initialization from the Xavier initialization scheme.

H.1 Fixed-Dim

While we highlighted the results with the Gaussian mixture model and classification settings with only 2 clusters/classes, we also conducted experiments with an increased number of clusters and classes, making the problem even more challenging. Tables 7-14 shows that both forward and reverse KL methods perform reasonably, with forward KL struggling more in challenging scenarios.

Next, we also consider harder tasks based on the Bayesian Neural Network (BNN) paradigm, where we consider nonlinear regression and classification setups with different activation functions: TANH and RELU for a 1-layered and 2-layered BNN. We provide the results of our experiments in Tables 7-14. The results indicate that forward KL approaches struggle a lot in such scenarios, often achieving performance comparable to random chance. On the contrary, we see that reverse KL-based amortization leads to performances often similar to dataset-specific optimization, thereby showing the superiority of our proposed method.

H.2 Variable-Dim

Our experiments on variable dimensional datasets can be evaluated for arbitrary feature cardinality, of which we show a few examples in Section 4. In this section, we provide results for additional dimensionality setups. In particular, we refer the readers to Tables 15-22, which contain experimental results w.r.t different dimensionalities (e.g. 50D setup), as well as different number of clusters and classes, respectively, for the GMM

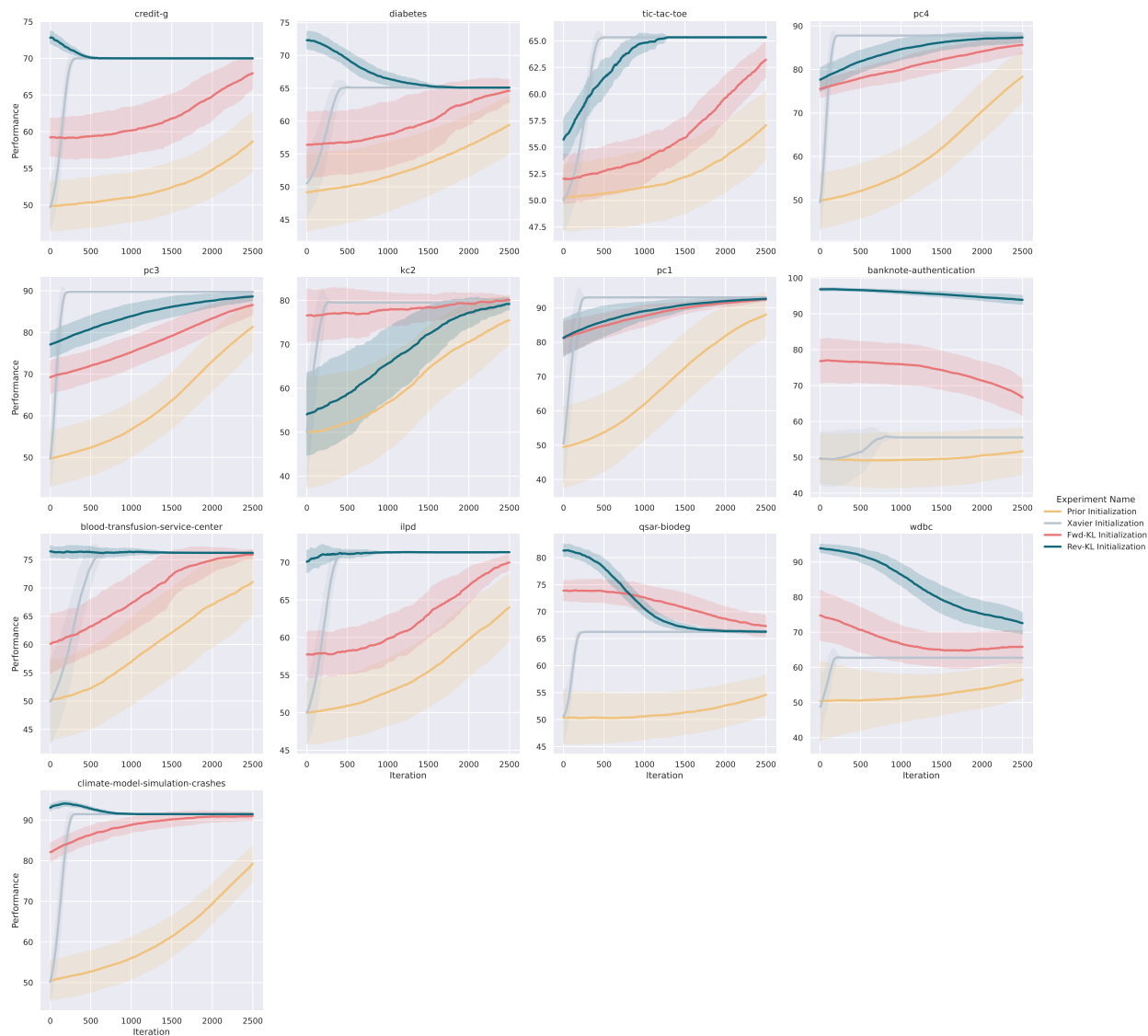


Figure 12: **Tabular Experiments | Linear Classification with Diagonal Gaussian:** For every classification dataset from the OpenML platform considered, we initialize the parameters of a linear classification-based probabilistic model with the amortized inference models which were trained with a diagonal Gaussian assumption. The parameters are then further trained with maximum-a-posteriori (MAP) estimate with gradient descent. Reverse and Forward KL denote initialization with the correspondingly trained amortized model. Prior refers to a MAP-based optimization baseline initialized from the prior $\mathcal{N}(0, I)$, whereas Xavier refers to initialization from the Xavier initialization scheme.

and LC setup. Throughout, we see that amortization leads to reasonable performance, and in particular, we see forward KL-based amortization starting to struggle in high-dimensional setups.

Again, to make the setup more challenging, we consider the Bayesian Neural Network (BNN) setup where we consider nonlinear regression and classification with different activation functions: TANH and RELU for a 1-layered and 2-layered BNN, but which can now be tested for an arbitrary number of input features. Our experiments are highlighted in Tables 15-22, for 1- and 2-layered BNN, among others. In such complex multi-modal and complicated setups, forward KL often performs comparable to random chance and thus does not lead to any good approximation of the true posterior distribution. On the other hand, our proposed method indeed leads to good predictive performance, often comparable to dataset-specific optimization routines.

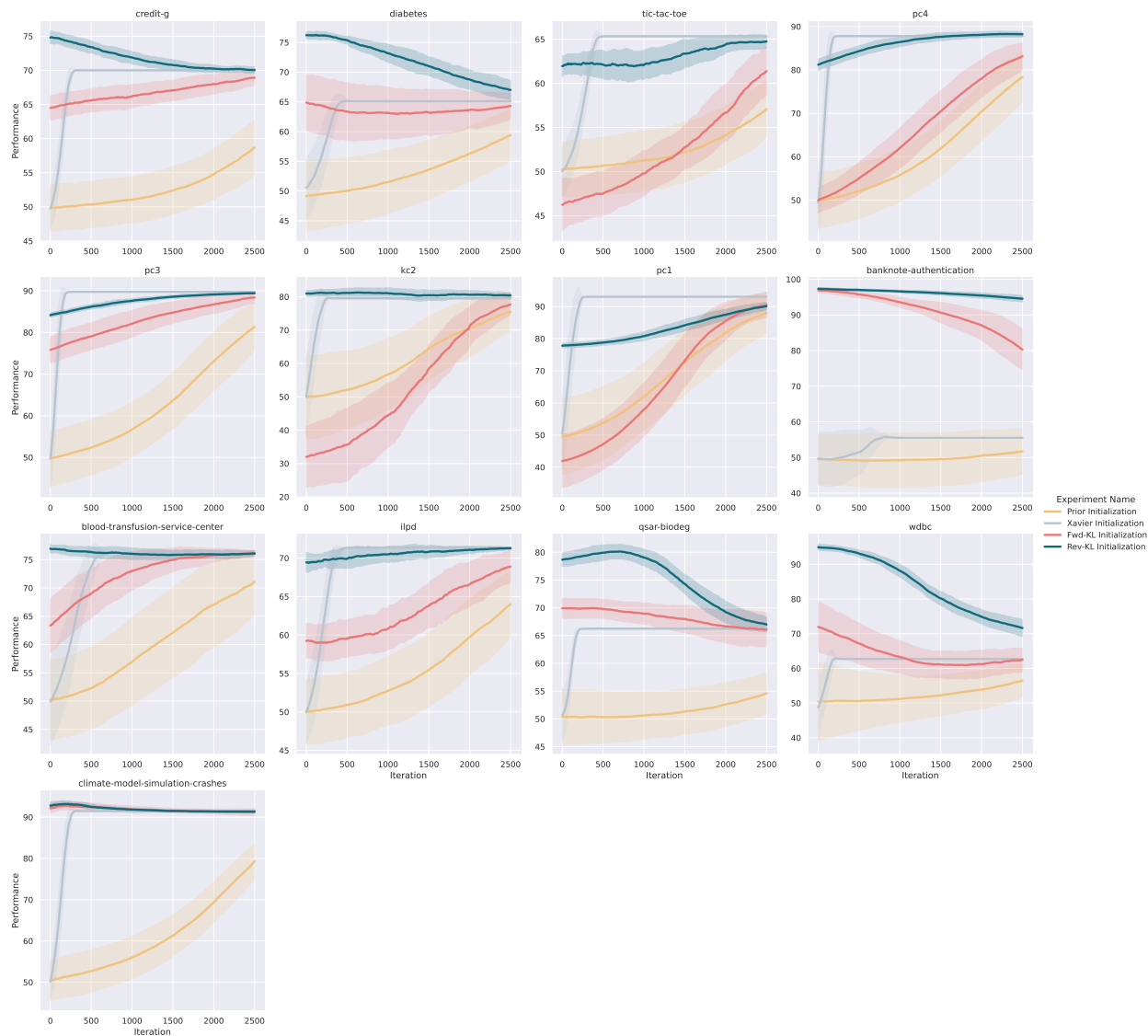


Figure 13: **Tabular Experiments | Linear Classification with Normalizing Flow:** For every classification dataset from the OpenML platform considered, we initialize the parameters of a linear classification-based probabilistic model with the amortized inference models which were trained with a normalizing flow-based model. The parameters are then further trained with maximum-a-posteriori (MAP) estimate with gradient descent. Reverse and Forward KL denote initialization with the correspondingly trained amortized model. Prior refers to a MAP-based optimization baseline initialized from the prior $\mathcal{N}(0, I)$, whereas Xavier refers to initialization from the Xavier initialization scheme.

H.3 Model Misspecification

As a representative of the results on model misspecification (Section 4), we highlighted training and evaluation of the amortization models with Transformer backbone on a subset of in-distribution and OoD data-generating functions (Table 3) to show superiority in generalization of reverse KL trained system vs. forward KL based ones on OoD data but also to highlight that training a misspecified amortization model on OoD datasets directly with our approach results in even better posterior predictive performance.

In addition to those experiments, we also conducted a broader range of experiments utilizing DeepSets as the backbone, various OoD data-generating functions for training and evaluation of the reverse KL system, and an additional nonlinear regression model with RELU activation function. For a comprehensive description

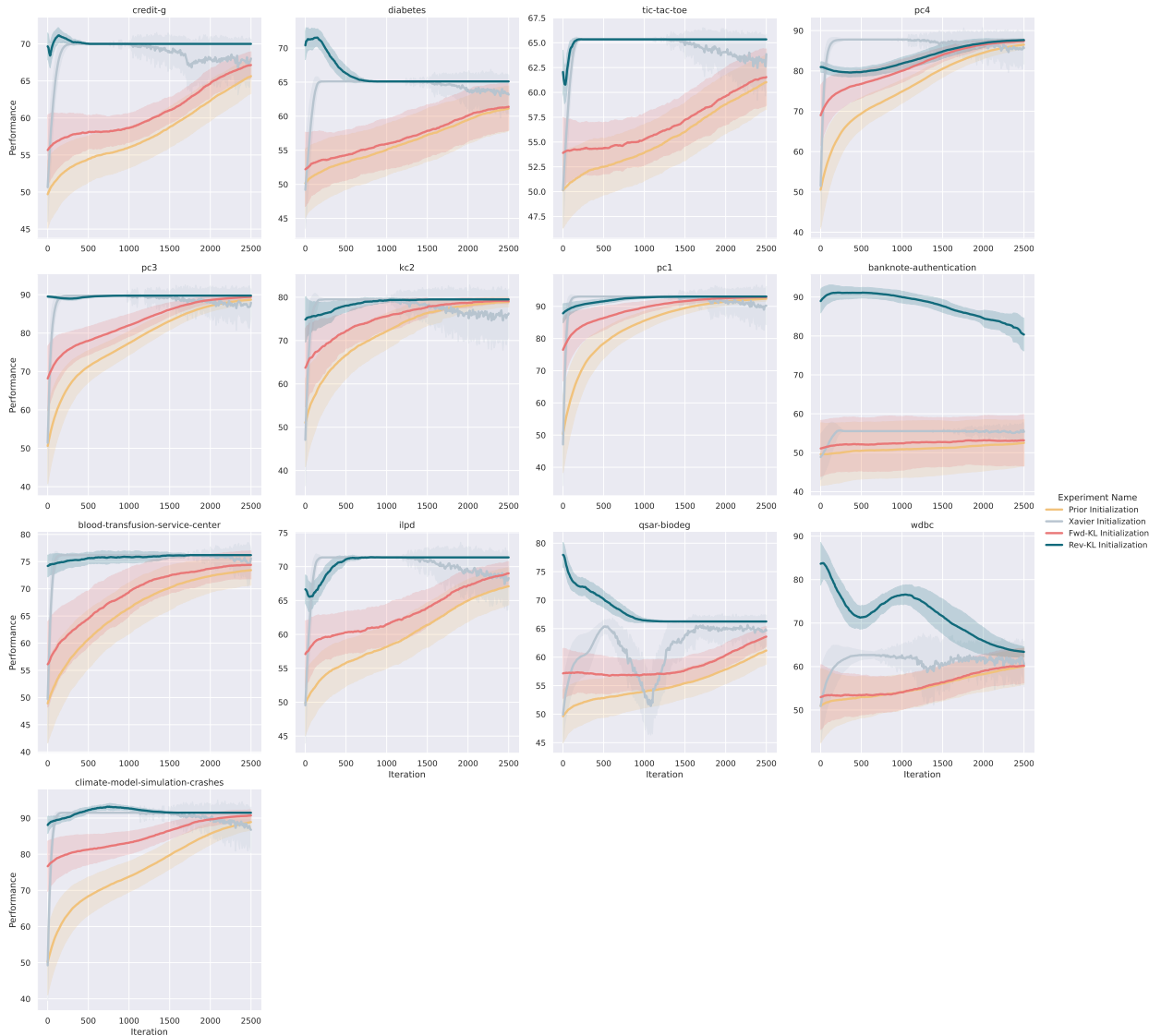


Figure 14: **Tabular Experiments | Nonlinear Classification with Diagonal Gaussian:** For every classification dataset from the OpenML platform considered, we initialize the parameters of a nonlinear classification-based probabilistic model with the amortized inference models which were trained with a diagonal Gaussian assumption. The parameters are then further trained with maximum-a-posteriori (MAP) estimate with gradient descent. Reverse and Forward KL denote initialization with the correspondingly trained amortized model. Prior refers to a MAP-based optimization baseline initialized from the prior $\mathcal{N}(0, I)$, whereas Xavier refers to initialization from the Xavier initialization scheme.

of these experiments and the complete setup, please refer to Section F.3. We considered two probabilistic models, including a linear regression model and a nonlinear regression models utilizing the TANH activation function. The detailed results for each model can be found in Tables 23 and 24.

In all experiments, reverse KL outperforms forward KL trained amortization models in in-distribution performance and excels in posterior prediction on OoD datasets. Although the significant difference in posterior prediction performance of forward vs. reverse KL in cases where the underlying model is nonlinear was already mentioned in previous experiments, here, reverse KL-trained models also excel in evaluations of posterior prediction for the linear regression model. Although only by a margin, in the case of approximating the posterior of the simpler linear regression model, a diagonal Gaussian-shaped posterior shows the best posterior prediction results when evaluated on OoD datasets from the nonlinear regression dataset generating

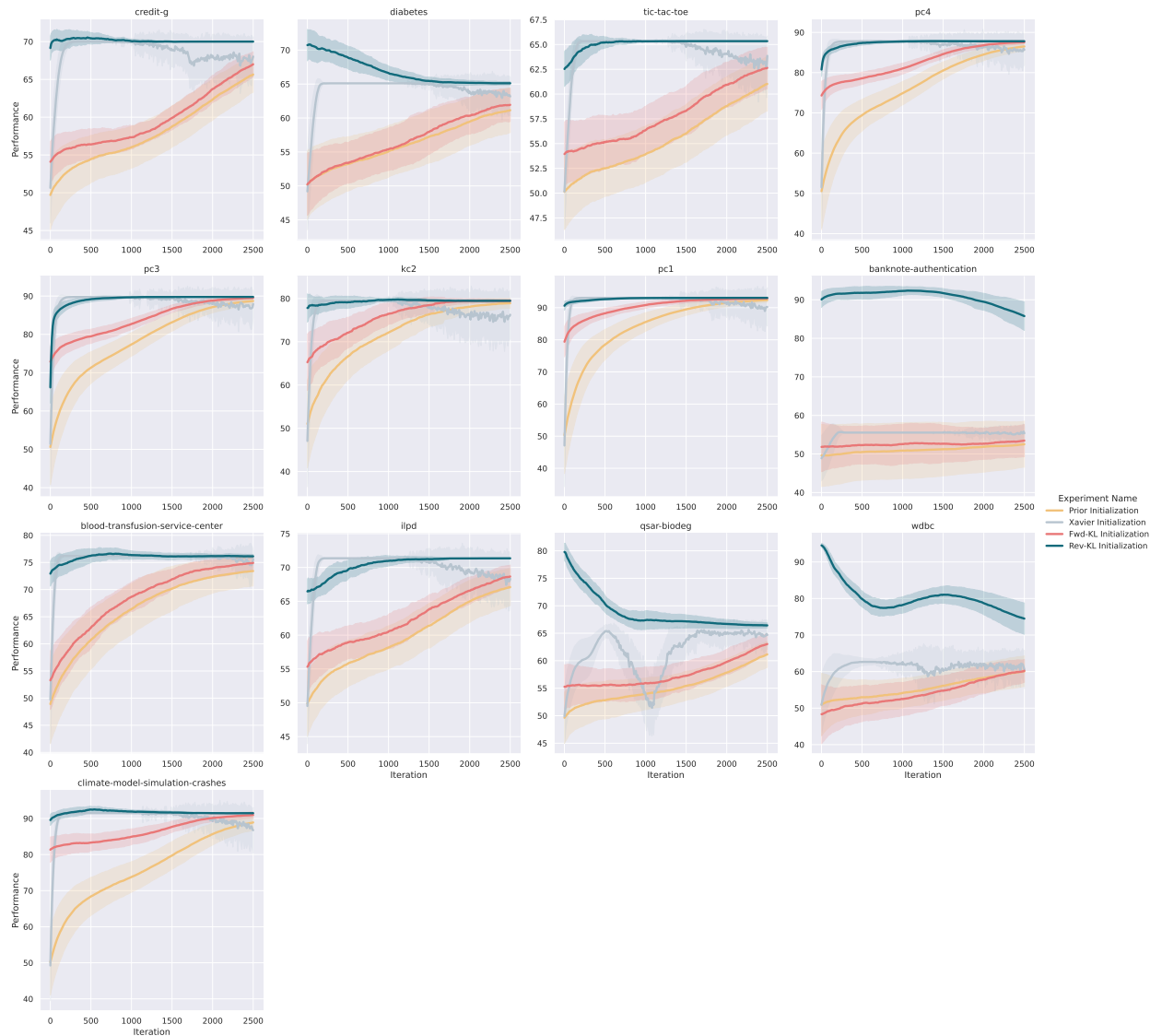


Figure 15: **Tabular Experiments | Nonlinear Classification with Normalizing Flow**: For every classification dataset from the OpenML platform considered, we initialize the parameters of a linear classification-based probabilistic model with the amortized inference models which were trained with a normalizing flow-based model. The parameters are then further trained with maximum-a-posteriori (MAP) estimate with gradient descent. Reverse and Forward KL denote initialization with the correspondingly trained amortized model. Prior refers to a MAP-based optimization baseline initialized from the prior $\mathcal{N}(0, I)$, whereas Xavier refers to initialization from the Xavier initialization scheme.

function. In almost all other experiments, the posterior prediction performance could be enhanced when we used the normalizing flow based posterior. A definitive conclusion cannot be drawn regarding the superiority of one backbone over the other, i.e. between DeepSets or Transformer. However, amortization models with DeepSets as the backbone tend towards better generalization regarding OoD datasets.

H.4 Tabular Experiments

As a case of extreme OoD generalization, we test our amortized models trained to handle variable feature dimensions on the suite of regression and classification problems that we filtered out from the OpenML platform, as outlined in Appendix G. We consider both linear and nonlinear probabilistic models to tackle the regression

and binary classification setups, which lead to predicting the parameters of a linear regression/classification model and a small nonlinear neural network based on RELU activation function. Further, we also perform the analysis with a diagonal Gaussian assumption and a normalizing flow-based amortization model trained with both a forward and reverse KL objective. We provide the results on the regression problems in (a) linear model with diagonal Gaussian assumption (Figure 8), (b) linear model with normalizing flow (Figure 9), (c) nonlinear model with diagonal Gaussian assumption (Figure 10), and (d) nonlinear model with normalizing flow (Figure 11). The results of the classification problems are shown in (a) linear model with diagonal Gaussian assumption (Figure 12), (b) linear model with normalizing flow (Figure 13), (c) nonlinear model with diagonal Gaussian assumption (Figure 14), and (d) nonlinear model with normalizing flow (Figure 15). Our experiments indicate that initializing with amortized models leads to better performance and training than models trained via maximum a-posteriori approach and initialized with the prior, i.e., $\mathcal{N}(0, I)$.

We do provide an additional baseline of initializing with XAVIER-INIT initialization, which often leads to faster convergence; however, as we consider the prior to be a unit normal, this is an unfair baseline as we assume the weights to be initialized from a different prior. We leave the work of computing Bayesian posteriors with different priors and testing an amortized Bayesian model with XAVIER-INIT prior for the future.