756 A MODEL MISSPECIFICATION

758 A.1 MIS-CALIBRATION VS MISSPECIFICATION 759

To further elucidate the distinction between posterior calibration and model misspecification, it is
 essential to highlight their respective scopes and the specific challenges they address.

762 Posterior calibration focuses on ensuring that the predicted posterior distributions accurately reflect the 763 true uncertainty in parameter estimates given the observations, under the assumption that the simulator 764 is well-specified. Methods such as those proposed by Falkiewicz et al. (2024); Delaunoy et al. (2022) address this by improving the alignment between the expected and actual coverage probabilities of the 765 posterior. These approaches generally assume that the simulator faithfully represents the generative 766 process of the observed data, enabling calibration to be evaluated and improved by leveraging 767 simulations. While important, these methods do not account for discrepancies between the simulator 768 and real-world data, which are precisely the scenarios we target in this work. 769

Model misspecification, on the other hand, arises when the simulator fails to capture the true generative process underlying the observed data. This results in systematic discrepancies that cannot be corrected solely by optimizing posterior calibration techniques. Misspecification introduces a gap between the simulated and real-world distributions, and this gap is only observable when real-world data is available. Unlike posterior calibration, addressing misspecification requires methods that can robustly leverage the simulator despite its inaccuracies, while incorporating real-world observations to mitigate the impact of the mismatch.

In our work, we explicitly focus on handling model misspecification. This distinction is reflected in the design of our approach and the evaluation scenarios we consider, such as Task E, where the simulated data diverges significantly from the real-world measurements. While posterior calibration methods may perform well in a well-specified context, they are not designed to cope with such gaps.
Instead, we prioritize creating predictive models that balance informativeness and robustness in the presence of misspecification, even if achieving perfect calibration remains an open and challenging problem.

784 A.2 COMPARISON BETWEEN MODEL MISSPECIFICATION DEFINITIONS

785
786 We provide a toy example to show how a simulator may be well-specified according to the standard definition of misspecification but still provide biased estimates of the target parameter when applied to real data.

Consider the following setting: a noisy sensor measures some physical quantity θ , producing measurements $\mathbf{x}_{o}^{1}, \ldots, \mathbf{x}_{o}^{n} \stackrel{\text{i.i.d.}}{\sim} \mathbb{P}^{\star}$, where $\mathbb{P}^{\star} := \mathcal{N}(\theta^{\star}, 1)$ is a normal distribution centered around the 'true' value θ^{\star} . Let { $\mathbb{P}_{\theta} : \theta \in \mathbb{R}$ } be a simulator of this process with $\mathbb{P}_{\theta} := \mathcal{N}(\mu, 1)$, where $\mu := \theta + \lambda$ and $\lambda > 0$ is a fixed scalar constant, which is a misspecification in the simulator that falsely accounts for a non-existing offset in the sensor that produced the real observations $\mathbf{x}_{o}^{1}, \ldots, \mathbf{x}_{o}^{n}$.

According to the standard definition of misspecification, the simulator is well specified, as setting $\theta \leftarrow \theta^* - \lambda$ yields $\mathbb{P}_{\theta} = \mathbb{P}^*$. However, the posterior estimates we obtain with this simulator are biased with respect to the true parameter θ^* .

To see this, let us compute the posterior under a Gaussian prior $\mathcal{N}(\theta^*, 1)$ over the parameter θ , centered on the true value θ^* . Taking advantage of the conjugate prior, the posterior $p(\theta \mid \mathbf{x}_o^1, \dots, \mathbf{x}_o^n)$

- 800
- 801 802

803

- 804
- 805
- 806
- 807
- 808

810	becomes	
811	$p(\theta \mid \mathbf{x}_1^1, \dots, \mathbf{x}_n^n)$	$\mathbf{x} \ p(\theta) p(\mathbf{x}^1, \dots, \mathbf{x}^n \mid \theta)$
812	$P(\circ 10, \cdots, 10) \circ$	n
813		$= p(\theta) \prod p(\mathbf{x}_{e}^{i} \mid \theta)$
814		$r(v) \prod_{i=1}^{r(v)} r(v)$
815		$1 (1) \frac{n}{2} 1 (1)$
017		$= \frac{1}{\sqrt{2-\epsilon}} \exp\left(-\frac{1}{2}(\theta - \theta^{\star})^2\right) \prod \frac{1}{\sqrt{2-\epsilon}} \exp\left(-\frac{1}{2}(\mathbf{x}_o^i - \mu)^2\right)$
010		$\sqrt{2\pi}$ (2) $\int_{i=1}^{2} \sqrt{2\pi}$ (2)
010 910		$\begin{pmatrix} 1 & 1 & 1 \end{pmatrix} \begin{pmatrix} n & 1 \end{pmatrix} \begin{pmatrix} n & 1 \end{pmatrix} \begin{pmatrix} n & 1 \end{pmatrix}$
820		$\propto \exp\left(-\frac{1}{2}(heta- heta^{\star})^2 - \frac{1}{2}\sum(\mathbf{x}_o^{\iota}-\mu)^2 ight)$
821		
822		$= \operatorname{ovn} \left(-\frac{1}{n} \left[\theta^2 + (\theta^\star)^2 - 2\theta \theta^\star + \sum_{i=1}^n (\mathbf{x}^i)^2 + n \mu^2 - 2\mu \sum_{i=1}^n \mathbf{x}^i \right] \right)$
823		$= \exp\left(-\frac{1}{2}\left[0^{-1} + (0^{-1}) - 200^{-1} + \sum_{i=1}^{n} (\mathbf{x}_{o})^{-1} + n\mu^{-1} - 2\mu \sum_{i=1}^{n} \mathbf{x}_{o}\right]\right)$
824		$\begin{pmatrix} 1 \\ n \\ n \end{pmatrix}$
825	(drop const. terms)	$\propto \exp\left(-\frac{1}{2}\left[\theta^2 - 2\theta\theta^{\star} + n\mu^2 - 2\mu\sum \mathbf{x}_o^i\right]\right)$
826		$\left(\begin{array}{c}2\\2\end{array}\right)$
827		$\begin{pmatrix} 1 \\ n \\$
828	$(\mu = \theta + \lambda)$	$= \exp\left(-\frac{1}{2}\left \theta^2 - 2\theta\theta^{\star} + n\theta^2 + n\lambda^2 + 2n\lambda\theta - 2\theta\sum \mathbf{x}_o^i - 2\lambda\sum \mathbf{x}_o^i\right \right)$
829		$\begin{pmatrix} 2 \\ i=1 \end{pmatrix}$
830		$\begin{pmatrix} 1 \begin{bmatrix} a^2 & aaa + a^2 + a & a & a \end{bmatrix}$
831	(drop const. terms)	$\propto \exp\left(-\frac{1}{2}\left \theta^{2}-2\theta\theta^{2}+n\theta^{2}+2n\lambda\theta-2\theta\sum_{i=1}^{n}\mathbf{x}_{o}^{i}\right \right)$
832		
833		$-\exp\left(-\frac{1}{2}\left[(n+1)\theta^2 - 2\theta(\theta^* - n\lambda + \sum_{i=1}^n \mathbf{x}^i)\right]\right)$
834		$= \exp\left(-2\left[\frac{(n+1)^{o}}{2}\left[\frac{2}{(n+1)^{o}}\right]\right]\right)$
835		$\begin{pmatrix} 1 \\ 1 \end{pmatrix}$
836		$= \exp\left(-\frac{1}{2(n+1)}\left[\theta^2 - 2\theta\left(\frac{1}{n+1}\right)\left(\theta^{\star} - n\lambda + \sum \mathbf{x}_{\alpha}^i\right)\right]\right)$
837		$\left(\begin{array}{cc}2(n+1)^{-1} \\ (n+1) \\$
838		$\begin{pmatrix} 1 & \lceil (1) & n \rceil^2 \end{pmatrix}$
839	(complete square)	$\propto \exp\left(-\frac{1}{2(-\lambda_{c})^{-1}}\left \theta - \left(\frac{1}{-\lambda_{c}}\right)\left(\theta^{\star} - n\lambda + \sum_{i} \mathbf{x}_{c}^{i}\right)\right \right),$
840	· · · · /	$\begin{bmatrix} 2(n+1)^{-1} \\ n+1 \end{bmatrix} \begin{bmatrix} n+1 \end{bmatrix}$
ö41	that is a normal distr	ibution $\mathcal{N}(\tau \ \gamma^2)$ with mean
042	and is, a normal abdi	

$$\tau = \left(\frac{1}{1+n}\right) \left(\theta^{\star} - n\lambda + \sum_{i=1}^{n} \mathbf{x}_{o}^{i}\right)$$

and variance $\gamma^2 = (n+1)^{-1}$. Thus, the posterior is biased, e.g., the posterior mean τ is a biased estimator of θ^* with $\mathbb{E}[\theta^* - \tau] = \theta^* - \lambda\left(\frac{n}{n+1}\right)$.

B THE ROPE ALGORITHM

Algorithm 1 Posterior Inference using Robust Neural Posterior Estimation (RoPE)

Input: Simulator $S(\theta, \varepsilon)$, prior distribution $p(\theta)$, calibration set $\mathcal{C} = \{(\mathbf{x}_o^i, \theta^i)\}_{i=1}^{N_c}$, test set $\mathcal{D} = \{\mathbf{x}_o^i\}_{i=1}^{N_o}$ **Output:** $\tilde{p}(\theta \mid \mathbf{x}_o) \forall \mathbf{x}_o^i \in \mathcal{D}$ **Step 1: Neural Posterior Estimation (NPE)** Train neural network \mathbf{h}_{ω} and conditional normalizing flow $p(\theta \mid \cdot)$ using NPE: $\tilde{p}, \omega^{\star} = \arg \max_{p, \omega} \mathbb{E}_{\substack{\theta \sim \pi(\theta) \\ \varepsilon \sim \mathcal{U}[0, 1]}} \left[\log p(\theta \mid \mathbf{h}_{\omega}(S(\theta, \epsilon))) \right]$ Step 2: Fine-tune sufficient statistics h_{ω^*} on the Calibration Set $\mathbf{g}_{\psi} := \operatorname{COPY}(\mathbf{h}_{\omega^{\star}})$ $\mathcal{C}_{train}, \mathcal{C}_{val} = \text{RandomSplit}(\mathcal{C}, \frac{1}{5})$ $best_{val} = \infty$ for N_{iter} do $\psi \leftarrow \psi - \alpha \nabla_{\psi} \left| \sum_{(\theta, \mathbf{x}_o) \in \mathcal{C}_{train}} |\mathbf{g}_{\psi}(\mathbf{x}_o) - \mathbb{E}_{\varepsilon} [\mathbf{h}_{\omega^{\star}}(S(\theta, \varepsilon))]|_2 \right|$ $\mathbf{cur}_{val} = \sum_{(\theta, \mathbf{x}_o) \in \mathcal{C}_{val}} |\mathbf{g}_{\psi}(\mathbf{x}_o) - \mathbb{E}_{\varepsilon}[\mathbf{h}_{\omega^{\star}}(S(\theta, \varepsilon))]|_2$ if $cur_{val} < best_{val}$ then $best_{val} = cur_{val}$ $\psi^{\star} = \psi$ end if end for Step 3: Generate Simulations for Test Set $(N_s = N_o)$ $\mathcal{S} = \{\mathbf{x}_s^j\}_{j=1}^{N_s},$ where $\mathbf{x}_s^j \sim S(\theta^j, \varepsilon) \quad \theta^j \sim \pi(\theta) \quad \varepsilon \sim \mathcal{U}[0, 1]$ Step 4: Entropic-regularized OT $C_{ij} = |f_{\omega^{\star}}(\mathbf{x}_s^j) - g_{\psi^{\star}}(\mathbf{x}_o^i)| \quad \forall i, j \in \{1, \dots, N_o\} \times \{1, \dots, N_s\}$ $P^{\star} = \arg\min_{\boldsymbol{P} \in \mathcal{B}_o} \langle \boldsymbol{P}, \boldsymbol{C} \, \rangle + \rho \, KL\left(\boldsymbol{P}^T \mathbf{1}_{N_o} \| \frac{\mathbf{1}_{N_s}}{N_s}\right) + \gamma \langle \boldsymbol{P}, \log \boldsymbol{P} \rangle$ **Step 5: Compute Posterior Distributions**

$$p(\theta | \mathbf{x}_o^i) := \sum_{j=1}^{N_s} P_{ij}^{\star} \tilde{p}\left(\theta \mid \mathbf{h}_{\omega^{\star}}(\mathbf{x}_s^j)\right)$$

Return $\tilde{p}(\theta | \mathbf{x}_o^i) \quad \forall \mathbf{x}_o^i \in \mathcal{D}$

918 C ROBUSTNESS TO PRIOR MISSPECIFICATION 919

968

969

970

971

In some practical applications of our algorithm, it is unlikely that the prior used to generate synthetic
 data will match the distribution of the target parameters in the real data. For this reason, we
 consider a semi-balanced formulation of OT, providing the flexibility to discard simulations with no
 corresponding real-world observations.

Prior misspecification on Task E. To evaluate the effect of a misspecified prior on RoPE and RoPE*, we perform an experiment that would resemble its use in real applications like the ones we outline in the introduction. In such settings—e.g., inferring cardiac parameters or chemical concentrations—the target parameters are limited to a range of validity, and a likely choice for the practitioner would be to select a uniform prior over this range.

To replicate this setting, we collect a new real-world dataset from the light tunnel (Task E) and train RoPE on synthetic data originating from a uniform prior, as we do for the results shown in Figure 2. However, we then apply RoPE to real data generated from a different (betabinomial) distribution over the target parameters. The results are shown in Figure 4, together with a visualization of the misspecified and true parameter distributions (prior A and B, respectively). We also show the learnt posterior distributions in Figure 5 for both RoPE and RoPE* ($\tau = 0.5$).

The results show that RoPE is relatively robust to prior misspecification prior. Furthermore, using an unbalanced OT formulation significantly improves the performance in this setting.

Prior misspecification on Task C. With this experiment we aim to better understand the role of τ when RoPE is applied with different levels of prior misspecification. We thus re-use the same setup as in Figure 2 but add prior misspecification as a mixture between the assumed prior and a much tighter uniform distribution. As the weight of the tighter uniform distribution increases, the prior gets more misspecified. The experimental setup follows closely the one in the well-specified case (see subsection 1.2), except calibration samples are drawn from the true prior (as this would be the case in a real-world application) and we compute the OT coupling for values of $\tau \in [0.1, 1]$.



Figure 5: Visualization of estimated posteriors. Corner plots of the posteriors estimated by RoPE in the prior-misspecification experiment from Fig. 1 above. We show, in different colors, the estimates for four observations sampled at random from the test set, for RoPE (left) and RoPE^{*} ($\tau = 0.5$) (right) formulation of the OT step, and a calibration set of size 50; the horizontal and vertical lines correspond to the ground-truth value of the parameters.





Figure 6: Out-of-distribution performance of RoPE and some baselines. We train RoPE and other baselines on the same light-tunnel data as in task E (training distribution), but apply it to test sets originating from a target distribution where the real-world images are flipped vertically. We compare the performance on test sets from both distributions, showing the LPP and ACAUC scores for each method. For comparison, in the right plot we show again the LPP curve (light gray, dotted) attained by RoPE under the training distribution. The performance of RoPE is barely affected as it cannot exploit any signal in the real images (\mathbf{x}_{o}) beyond what is encoded in the simulator, and the simulator output (\mathbf{x}_s) is invariant to the transformation we consider. Because NPE is not trained on real observations, its performance, although poor, also remains virtually unchanged. On the other hand, the performance of MLP and J-NPE drops in the target distribution, as these methods are not limited in what information they can exploit from the real observations on which they are trained, potentially learning shortcuts that are not present in the target distribution. This results demonstrate that if the simulator embeds the right invariances, our modeling assumption $\mathbf{x}_o \perp \theta \mid \mathbf{x}_s$ can be favorable to out-of-distribution generalization.

Ε **OPTIMAL TRANSPORT COUPLING AS A JOINT DISTRIBUTION**

With our conditional independence assumption, the problem of modeling $p(\mathbf{x}_o \mid \theta)$ reduces to modeling $p(\mathbf{x}_o \mid \mathbf{x}_s)$ instead. If we assume the prior well-specified, this task is equivalent to modeling $p(\mathbf{x}_o, \mathbf{x}_s)$ under the constraint that the corresponding marginal $p(\mathbf{x}_s) = \int p(\mathbf{x}_s, \mathbf{x}_o) d\mathbf{x}_o$ equals $\int p(\theta)p(\mathbf{x}_s \mid \theta)d\theta$. By construction, the OT coupling, π^* , respects the constraint on the marginals, $\int \pi^*(\mathbf{x}_s, \mathbf{x}_o) d\mathbf{x}_o = p(\mathbf{x}_s)$ and $\int \pi^*(\mathbf{x}_s, \mathbf{x}_o) d\mathbf{x}_s = p(\mathbf{x}_o)$, and the exact instantiation π^* depends also on the chosen cost function which can always be defined to yield any given conditional $p(\mathbf{x}_o \mid \mathbf{x}_s)$ that respects the constraint $\int p(\mathbf{x}_o \mid \mathbf{x}_s) p(\mathbf{x}_s) d\mathbf{x}_s = p(\mathbf{x}_o)$. π^* can thus model the "right" posterior, provided the right cost function is used. In the case, where the prior cannot be trusted, we suggest to use $\tau < 1$ and relax the OT formulation. In this case, we only enforce that all elements of $p(x_o)$ are matched to a subset of the elements of $p(x_s)$. This implicitly assumes that the assumed prior $p(\theta)$ is overly conservative and covers $p^{\star}(\theta)$. We believe this is a reasonable assumption as it is often easy to derive physical bounds for the parameter values and use a uniform distribution.

1026 F **SELF-CALIBRATION PROPERTY** 1027

1030 1031

1037 1038 1039

104

1052

1028 We say RoPE is self-calibrating because, by design, the posterior distribution marginalized over observations tends to the prior as the number of simulation increases, that is, 1029

$$\int_{\mathcal{X}} \tilde{p}(\theta \mid \mathbf{x}_o) p(\mathbf{x}_o) d\mathbf{x}_o = p(\theta).$$
(7)

1032 This property is also called marginal calibration, and is a necessary condition for a posterior estimation 1033 method to be calibrated. Considering NPE, $\tilde{p}(\theta \mid \mathbf{x}_s)$, is marginally calibrated and observations 1034 \mathbf{x}_o are generated from the assumed prior, that is sampled from an unknown distribution $p(\mathbf{x}_o) =$ 1035 $\int p(\mathbf{x}_o \mid \theta) p(\theta)$, we can show RoPE is marginally calibrated. Indeed, considering the Monte-Carlo 1036 approximation of the marginalized posterior distribution over the test set $\mathcal{D}_o := \{\mathbf{x}_o^i\}_{i=1}^{N_o}$, we have,

$$\int_{\mathcal{X}} \tilde{p}(\theta \mid \mathbf{x}_o) p(\mathbf{x}_o) d\mathbf{x}_o = \mathbb{E}_{p(\mathbf{x}_o)}[\tilde{p}(\theta \mid \mathbf{x}_o)]$$
(8)

1040
1041
$$\approx \frac{1}{N_o} \sum_{i=1}^{N_o} \tilde{p}(\theta \mid \mathbf{x}_o^i) \tag{9}$$
1042

1042
1043
1044
1045

$$= \frac{1}{N_o} \sum_{i=1}^{N_o} \sum_{j=1}^{N_s} N_o P_{ij}^{\star} \tilde{p}(\theta \mid \mathbf{x}_s^j)$$
(10)

1046
1047
1048
1049
$$=\sum_{j=1}^{N_s} \left[\sum_{i=1}^{N_o} P_{ij}^{\star}\right] \tilde{p}(\theta \mid \mathbf{x}_s^j) \tag{11}$$

1049
1050
$$= \frac{1}{N_s} \sum_{j=1}^{N_s} \tilde{p}(\theta \mid \mathbf{x}_s^j)$$
(12)
1051

$$\approx p(\theta),$$
 (13)

where we use the definition of the transport matrix to get $\sum_{i=1}^{N_o} P_{ij}^{\star} = \frac{1}{N_s}$. The last approximation tends to be exact as the number of simulations increases, if the NPE is marginally calibrated. 1053 1054 1055

1056 G LEARNING MINIMAL SUFFICIENT STATISTICS WITH NEURAL POSTERIOR 1057 **ESTIMATION** 1058

We now discuss why NPE may learn a minimal sufficient statistic under perfect training. First, under 1059 a sufficiently large validation set, NPE's objective function is only optimal on the validation set if NPE models the true posterior as defined implicitly by the prior $p(\theta)$ and the likelihood corresponding 1061 to the simulator S. This consistency has been proven in (Papamakarios & Murray, 2016) and is 1062 the motivation to use such an objective when estimating density. Second, some normalizing flows, 1063 such as autoregressive UMNN flows (Wehenkel & Louppe, 2019), are universal approximators of 1064 continuous densities. In addition, neural networks are also universal function approximators. As such, we can claim that it is always possible to parameterize the NCDE $p_{\theta}(\theta \mid \mathbf{h}_{\omega}(\mathbf{x}))$ such that the class of functions its parameters represent contains the true posterior. We directly observe that \mathbf{x} is 1067 only used by the NCDE through $\mathbf{h}_{\omega}(\mathbf{x})$. Thus, under perfect training $p_{\theta^{\star}}(\theta \mid \mathbf{h}_{\omega^{\star}}(\mathbf{x})) = p(\theta \mid \mathbf{x})$ 1068 and $\mathbf{h}_{\omega^{\star}}(\mathbf{x})$ is a sufficient statistic for θ given \mathbf{x} under the simulator's model.

1069 Without additional constraints, we cannot claim anything about the minimality of $h_{\omega^*}(\mathbf{x})$. Neverthe-1070 less, we can enforce the neural network $\mathbf{h}_{\omega^{\star}}(\mathbf{x})$ to have an information bottleneck and thus reduce 1071 the information carried. In practice, we choose the output dimension of $\mathbf{h}_{\omega^{\star}}(\mathbf{x})$ so that the NCDE 1072 achieves optimal performance on the test set. Because in the context of SBI we can generate as many (simulated) samples as needed, we can obtain estimators that closely approach the simulation's 1074 posterior and a minimal sufficient statistic.

1075

1077

1078

¹⁰⁸⁰ H COMPUTATIONAL COST OF ROPE

1082 Running NPE is broadly recognized as having a low computational cost: once the upfront training is complete, the cost of inverting the normalizing flow to sample from the posterior during inference becomes negligible as the number of test observations increases. This makes NPE more efficient 1084 than methods like Approximate Bayesian Computation or Markov Chain Monte Carlo (when the simulator allows likelihood evaluation). RoPE introduces additional computational costs on top of 1086 running NPE: (1) the OT coupling computation, i.e., solving Equation 2, and (2) obtaining samples 1087 from the estimated posterior distributions, to compute the posterior estimate defined in Equation 5. 1088 The computational cost of solving the transport problem with the Sinkhorn algorithm (Cuturi, 2013) 1089 is quadratic in the number of real-world observations. The sampling step has a negligible cost as it 1090 directly sub-samples from the set of points generated with NPE. 1091

In our experiments, solving the OT optimization for 2000 test examples takes less than a minute on an M1 MacBook Pro. Sampling from the mixture of posterior distributions involves caching 10,000 samples for each simulation and generating 5,000 samples by sub-sampling from the mixture using the OT coupling matrix. This caching process takes under three minutes, and is comparable to the cost of running NPE alone.

Extending RoPE to handle larger test sets or an online setting (processing test examples one at a time) is outside the scope of this work. Nevertheless, mehtods like Neural OT (e.g., (Makkuva et al., 2020)) and online Sinkhorn (Mensch & Peyré, 2020) should provide good solutions to make RoPE fully amortized.

1100 1101

1127

I EXPERIMENTAL SETUP

In this section, we provide more details on our experiments. For completeness, we provide details on the neural architectures and training hyperparameters. However, we encourage the reader interested in reproducing our experiments to examine our code directly (a link to the code will be made available in the public version of the paper).

For all methods training on calibration set we keep always keep 20% of the calibration to monitor
 validation performance and we select the best model based on this metric.

For the MLP we use the same architecture as the NSE for all our experiments and optimize its parameters on the calibration set with Adam and a learning rate equal to 0.0003, we select the best model based on the LPP attributed to the validation subset of the calibration set.

Computing the SBI baseline. We take the ground-truth labels $\{(\theta^i\}_{i=1}^N \text{ from the test set } \{\theta^i, \mathbf{x}_o^i\}_{i=1}^N \}$ on which we compute all the metrics for Figure 2; for each label θ^i , we simulate a synthetic observation $\mathbf{x}_s^i := S(\theta^i)$, collecting them into a "synthetic" test set $\{(\theta^i, \mathbf{x}_s^i)\}_{i=1}^N$; then, we apply to it the NSE+NPE pipeline (simulated posterior in Figure 1) right) to obtain the posterior estimates which we then evaluate. In this way, the baseline represents the performance we would hope to achieve if there was no misspecification and the simulator perfectly replicated the real observations (up to the stochasticity of the simulator itself).

1120 I.1 TASK A: CS & TASK B: SIR

Task A (synthetic): CS. We reproduce the cancer and stromal cell development benchmark from Ward et al. (2022). The simulator emulates the development of cancer and stromal cells in a 2D environment as a function of three Poisson rate parameters (λ_c , λ_p , λ_d). The observations are vectors composed of the number of cancer and stromal cells and the mean and maximum distance between stromal cells and their nearest cancer cell. Synthetic misspecification is introduced by removing cancer cells that are too close to their generating parent.

Task B (synthetic): SIR. We also use the stochastic epidemic model from Ward et al. (2022), which describes epidemic dynamics through the infection rate β and recovery rate γ . Each observation is a vector composed of the mean, median, and maximum number of infections, the day of occurrence of the maximum number of infections, the day at which half the total number of infections was reached, and the mean auto-correlation (lag 1) of the infections. Misspecification is a delay in weekend infection counts, of which 5% are added to the count of the following Monday. We refer the reader to Ward et al. (2022) for more details about the simulator and prior distribution. We use the exact same setting as theirs.

1136 1137

1138 NEURAL ARCHITECTURE & TRAINING HYPERPARAMETERS

For all methods we use the same backbone MLP as the NSE with ReLU activations and layers composed of [4K, 16K, 16K, 12K, 3K] neurons, where K is the dimensionality of θ . The NF is a 1-step UMNN-MAF (Wehenkel & Louppe, 2019) with [100, 100, 100] neurons for both the autoregressive conditioner and normalizer. For NNPE, we train the UMNN-MAF on simulations poluted by Spike and Slab errors. We train models with Adam and a learning rate equal to 0.0005 and all other parameters set to default. We optimize the SBI model for 10^6 gradient steps and select the best model on random validation sets containing 10^5 simulations.

1146 I.2 TASK C: PENDULUM

1148 DESCRIPTION

1149 т

The first task is inspired from the damped pendulum benchmark commonly used to assess hybrid learning algorithms. Given a 2D physical parameter $\theta := [\omega_0, A]$, where $\omega_0 \in \mathbb{R}^+$ denotes the fundamental frequency and $A \in \mathbb{R}^+$ the amplitude of a friction-less pendulum, the simulator generates the horizontal position of the pendulum at 200 discrete times during uniformly sampled in a 10 seconds interval as

$$\mathbf{x}_s := [\theta(t=0), \dots, \theta(t=10s)] \in \mathbb{R}^{200}$$

where $\theta(t) = A \cos(\omega_0 t + \varphi) \quad \varphi \sim \mathbb{U}(-\pi, \pi).$ (14)

1157 The relationship between the parameters and the simulation is thus stochastic as φ accounts for an 1158 unknown phase shift when the measurements start. We generate real-world observations synthetically 1159 by replacing $\theta(t)$ from (14) by

$$\tilde{\theta}(t) = e^{\alpha t} A \cos(\omega_0 t + \varphi) \quad \varphi \sim \mathbb{U}[-\pi, \pi] \quad \alpha \sim \mathbb{U}[0, 1]$$

where α represents the effect of friction. We also add Gaussian noise on both simulated and real-world data to represent the inaccuracy of a sensor measuring the pendulum's position. The prior distribution is a product of uniform distribution, $p(\theta := [\omega_0, A]) = \mathcal{U}[0, 3] \times \mathcal{U}[0.5, 10]$.

1164 1165

1170

1160

1155 1156

1166 NEURAL ARCHITECTURE & TRAINING HYPERPARAMETERS

Neural Posterior Estimator. The NSE is a 1D convolutional neural network, with the architecture described in Algorithm 2. The NCDE is a one-step discrete normalizing flow with an autoregressive

Alg	orithm 2 Convolutional Neural Network for T	asks	A and D.	
1:	Conv1d(1, 16, 3, 1, dilation = 2, padding =	12:	ReLU()	
	1)	13:	Conv1d(128, 128, 3, 2, dilation	=
2:	ReLU()		2, padding = 1)	
3:	Conv1d(16, 64, 3, 2, dilation = 2, padding =	14:	ReLU()	
	1)	15:	AvgPool1d(3,1)	
4:	ReLU()	16:	Conv1d(128, 128, 3, 1, dilation	=
5:	AvgPool1d(3,1)		2, padding = 1)	
6:	Conv1d(64, 128, 3, 1, dilation = 2, padding =	17:	ReLU()	
	1)	18:	Flatten()	
7:	ReLU()	19:	Linear(2048, 512)	
8:	Conv1d(128, 128, 3, 2, dilation) =	20:	ReLU()	
	2, padding = 1)	21:	Linear(512, 128)	
9:	ReLU()	22:	ReLU()	
10:	AvgPool1d(3,1)	23:	Linear(128, 32)	
11:	Conv1d(128, 128, 3, 1, dilation) =	24:	ReLU()	
	2, padding = 1)	25:	Linear(32, 10)	

1188 1189	Algorithm 2 UNat1D Architecture	Algorithm 4 Plogh1D(in shannals
1190		out channels)
1191	1: UnetID:	1: Convid(in abannals out abannals
1192	2: Elicoueritz : 3: Block(in channels — 1 out channels —	kernel size=3 padding=1)
1193	$5.$ Diver (in_enamers = 1, out_enamers = 64)	2. ReLU()
1194	4: Block(in channels = 64 , out channels =	3: Conv1d(out channels, out channels,
1195	128)	kernel_size=3, padding=1)
1196	5: $Block(in_channels = 128, out_channels =$	4: ReLU()
1197	256)	
1198	6: Block(in_channels = 256 , out_channels =	
1199	512)	Algorithm 5 2D Convolutional Neural Net-
1200	7: Block(in_channels = 512 , out_channels =	work
1201	$\frac{1024}{\text{MaxPacl1d}(2)}$	1: Conv2d(3 64 3 2 dilation -1) ReLU()
1202	$\begin{array}{ccc} 0 & 0 \\ 0 & \mathbf{D} \\ 0 & 0 \end{array}$	$2: \text{Conv}2d(64 \ 128 \ 3 \ 2 \ \text{dilation=1}), \text{ ReEO()}$
1203	10: ConvTranspose1d(1024 $+$	ReLU()
1204	5 512 2 stride = 2)	3: MaxPool2d(3)
1205	11: Block(in channels $=$	4: Conv2d(128, 128, 3, 2, dilation=1),
1206	1024 , out_channels = 512)	ReLU()
1207	12: ConvTranspose1d $(512, 256, 2, stride = 2)$	5: Conv2d(128, 64, 1, 1, dilation=1),
1208	13: $Block(in_channels = 512, out_channels =$	ReLU()
1209	256)	6: Conv2d(64, 3, 1, 1, dilation=1), ReLU()
1210	14: ConvTranspose1d(256 , 128 , 2 , stride = 2)	7: Flatten() $2 - L_{1}^{1} + (27 - 100) - D - L H(0)$
1211	15: Block(in_channels = 256 , out_channels =	8: Linear(27 , 100), ReLU()
1212	128)	9: Linear(100, 20)
1213	16: Conviranspose $Id(128, 64, 2, stride = 2)$ 17: Block(in channels = 128 out channels =	
214	$17.$ Block (m_channels = 128, out_channels = 64)	
215	18: ConvTranspose1d(64, 1, 2, stride = 2)	
216	19: Block(in channels = 64 , out channels =	
217	1)	
218	20: $Conv1d(64, 1, 1)$	
1219		
1220		
1221		
1222	conditioner and a UMINN (Wehenkel & Louppe, 2	(2019) as the normalizer. The autoregressive
1223	conditioner is a MADE with ReLU activation and 3 lay	net with 2 layers of 100 neurons with Pal U
224	activations. For training the NDE, we use a batch size	of 100 and a learning factor equal to 1e 4 NPE
1225	is trained until convergence. Other parameters are set	to default values and should marginally impact
226	the NPE obtained	to default values and should marginarry impact
1227		1 1 11 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
1228	ROPE NSE. We have selected the best NPE based	1 on the validation set with 10000 examples
1229	generated with the simulator. The NPE is fixed to or	ne best-of-all model. We fine-tune the NCDE
1230	with a learning rate equal to 1e-5 for 5000 gradient s	steps on 80% the full calibration set. We use a
1231	1-sample Monte Carlo estimate of the expectation in	(O).
1232	J-NPE. To train J-NPE, we simply randomly use a bat	tch composed of 50% of simulated pairs (θ, \mathbf{x}_s)
1233	and of 50% (θ , \mathbf{x}_o) from the calibration set. We use t	he same architecture and hyper-parameters as
1234	the SBI NPE. The best model is selected based on the	best training set performance. We do 50 epochs
1235	with 50000 simulated examples for each epoch. The	batch size is 100.
1236	HVAE. For the HVAE, we re-use the NPE model as the	he physics encoder and replace the decoder with
1237	a deterministic version of the simulator, thus removing	the Gaussian noise on a random phase shift. In
1238	addition, we follow the approach of Takeishi & Kalous	sis (2021) and have 1) a real-world encoder that

¹²⁴¹ 5D latent vector \mathbf{z}_a . The reality-to-physics and physics-to-reality also have the same architectures and are two conditional 1D U-Net with neural network architecture described in Algorithm [3].

1239

1240

maps \mathbf{x}_o to \mathbf{z}_a , 2) a reality-to-physics encoder, and 3) a physics-to-reality decoder. The real-world

encoder has the same architecture as the NSE of the NPE and outputs the mean and log-variance of a

1242 To train the HVAE, we freeze the parameters of the NPE and optimizes the ELBO as well as a 1243 calibration loss that evaluates the likelihood assigned to the true physical parameters. All distributions 1244 are parameterized by Gaussian with mean and log-variance predicted by the neural networks. We 1245 do not use any additional losses as we expect constraining NPE and using the calibration set should 1246 already provide the necessary support to use the physics in a meaningful way. The HVAE is trained on the 2000 test examples as it is the only real-world data, calibration set aside, that we have access 1247 to. We use a batch size equal to 100 and a learning rate equal to 1e-3. We believe obtaining a better 1248 HVAE is possible. However, we emphasize the complexity of setting up a good HVAE for the only 1249 purpose of statistical inference over parameters. 1250

1251 1252 DATASETS

For this task, we can generate samples (θ, \mathbf{x}_s) on the fly to train the NPE. The calibration and test sets are also generated randomly by sampling from the prior distribution and using the damped pendulum simulator.

- 1257 I.3 TASK D: HEMODYNAMICS
- 1258 1259 DESCRIPTION

1260 Inspired by Wehenkel et al. (2023), we define the task of inferring important cardiovascular parameters 1261 from normalized arterial pressure waveforms measured at the radial artery. The simulator uses many 1262 physiological parameters that modulates the heart function, physical properties of the 116 main 1263 arterial segments, and behavior of the vascular beds. Our inference concerns two parameters of the heart function, $\theta := [SV, LVET]$, the stroke volume (SV) is the amount pumped out from the 1264 1265 left ventricle over the heart beat modeled, and the left ventricular ejection time (LVET) is the time interval between opening and closure of the aortic valve. Other parameters, such as the heart rate 1266 or arteries' stiffness, are considered as nuisance effects and are randomly sampled from a realistic 1267 population distribution. An additional source of randomness is added by modeling measurement 1268 errors with a white Gaussian noise and randomizing the starting recording time with respect to the 1269 cardiac cycle. The simulator produces 8-second timeseries $\mathbf{x}_t \in \mathbb{R}^{1000}$ sampled at 125Hz. As 1270 synthetic misspecification, the simulator assumes all arteries have the same length over the population 1271 considered, whereas "real-world" data are artificially generated by also varying the length of arteries 1272 and account for the effect of human's height. The simulator is based on the openBF PDE solver (Melis, 1273 (2017) specialized for hemodynamics, which is not differentiable and takes approximately one minute 1274 to simulate one sample on a standard CPU. This synthetic tasks represent a common scenario in 1275 which a simulator, although faithful to the effect of certain parameters, misses additional degrees of freedom that exists for the real-world data. 1276

1277 1278

1279

NEURAL ARCHITECTURE & TRAINING HYPERPARAMETERS

1280 Algorithm 6 CNN Architecture for Task C. 1281 1: Conv1d(1, 16, 3, 1, dilation=2, padding=1), ReLU() 1282 2: Conv1d(16, 64, 3, 2, dilation=2, padding=1), ReLU() 1283 3: AvgPool1d(4, 2) 1284 4: Conv1d(64, 128, 3, 1, dilation=2, padding=1), ReLU() 1285 5: Conv1d(128, 128, 3, 2, dilation=2, padding=1), ReLU() 1286 6: AvgPool1d(4, 2) 1287 7: Conv1d(128, 128, 3, 1, dilation=2, padding=1), ReLU() 8: Conv1d(128, 128, 3, 2, dilation=2, padding=1), ReLU() 9: AvgPool1d(4, 1) 10: Conv1d(128, 128, 3, 1, dilation=2, padding=1), ReLU() 1290 11: Flatten() 1291 12: Linear(1024, 512), ReLU() 13: Linear(512, 128), ReLU() 1293 14: Linear(128, 32), ReLU() 1294 15: Linear(32, 5) 1295

Neural Posterior Estimator. The NSE is the 1D convolutional neural network described in Algorithm 6. The NCDE is a 5-step discrete normalizing flow with an autoregressive conditioner and affine normalizers. Each of the 5 autoregressive conditioners is a MADE with ReLU activations and 4 layers of 300 neurons that output 4 dimensional vectors used to parameterize the affine transformations. For training the NPE, we use a batch size of 100 and a learning factor equal to 5e-4. NPE is trained until convergence. Other parameters are set to default values and should marginally impact the NPE obtained.

RoPE NSE. We have selected the best NPE based on the validation set with 2000 examples generated with the simulator. The NPE is fixed to one best-of-all model. We fine-tune the NCDE with a learning rate equal to 1e-5 for 2000 gradient steps on 80% of calibration set. We use a 1-sample Monte Carlo estimate of the expectation in (6).

J-NPE. To train J-NPE, we simply randomly use a batch composed of 50% of simulated pairs (θ, \mathbf{x}_s) and of 50% (θ, \mathbf{x}_o) from the calibration set. We use the same architecture and hyper-parameters as the SBI NPE. The best model is selected based on the best training set performance. We do 50 epochs with 6000 simulated examples for each epoch. The batch size is 100.

HVAE. There is no HVAE for this experiment as the simulator is non-differentiable.

- 1312 1313
- 1314 DATASETS

For this task, we cannot generate samples (θ, \mathbf{x}_s) on the fly to train the NPE. For the purpose of this experiment, we have generated 10000 simulations and real-world observations. Our fine-tuning strategy approximates (6) by finding the simulations with the closest parameter value.

- 1319 I.4 TASK E: LIGHT TUNNEL
- 1320 1321 DESCRIPTION

1322 We use one of the light-tunnel datasets from the causal chamber project (Gamella et al., 2024) 1323 causalchamber.org). In particular, we use the data from the ap_1.8_iso_500.0_ss_0 1324 005 experiment in the lt_camera_v1 dataset. The light tunnel is an elongated chamber with a 1325 controllable light source at one end, two linear polarizers mounted on rotating frames, and a camera 1326 that takes images of the light source through the polarizers. We refer the reader to Gamella et al. (2024) Figure 2) for a complete schematic. Our task consists of predicting the color setting of the light source 1327 $((R,G,B) \in [0,255]^3)$ and the dimming effect of the linear polarizers $\alpha \in [0,1]$ from the captured 1328 images. As a misspecified simulator of the image-generating process, we adopt the simple model 1329 described in Gamella et al. (2024, Model F1, Appendix D). A Python implementation is available 1330 through the causalchamber package (models.model_f1); visit causalchamber.org for 1331 more details. As input, the simulator takes the parameters $\theta := [R, G, B, \alpha]$ and produces an image 1332 consisting of a hexagon roughly the size of the light source, with an RGB color vector equal to 1333 $[\alpha R, \alpha G, \alpha B]$. The factor $\alpha := \cos^2(\theta_1 - \theta_2)$, where θ_1, θ_2 denote the angles of the two polarizers, 1334 corresponds to Malus' law (e.g., Collett, 2005), which models the dimming effect of the polarizers as 1335 a function of their relative angle. Besides the obvious misspecification with respect to image realism 1336 (see Figure 2), the model ignores other important physical aspects, such as the spectral response of 1337 the camera sensor or the non-uniform effect of the polarizers on the different colors—more details 1338 can be found in Gamella et al. (2024) Appendix D.IV.2.2). The prior is uniform over colors and polarizer angles, which leads to a non-uniform prior over the dimming effect α . 1339

1340

1342

1341 NEURAL ARCHITECTURE & TRAINING HYPERPARAMETERS

Neural Posterior Estimator. The NSE is the 2D convolutional neural network described by Algorithm 5.

The NCDE is also a one-step discrete normalizing flow with an autoregressive conditioner and a UMNN (Wehenkel & Louppe, 2019) as the normalizer. The autoregressive conditioner is a MADE with ReLU activation and 3 layers of 500 neurons that outputs a 10 dimensional vector to the UMNN. The UMNN has an integrand net with 4 layers of 150 neurons with ReLU activations. For training the NPE, we use a batch size of 100 and a learning factor equal to 5e-4. NPE is trained until convergence. Other parameters are set to default values and should marginally impact the NPE obtained.

Alg	orithm 7 2D UNet	Algorithm 8 Block2D(in_channels, out channels)
1:	Encoder2D: Block2D(in channels=3 out channels=64)	1: Conv2d(in channels out channels
2. 3.	Block2D(in_channels=64_out_channels=128)	kernel size=3 nadding=1 bias=False)
5.	Dioek2D(III_enalitieis=04, out_enalitieis=120)	2: BatchNorm2d(num features=out channels)
4:	Block2D(in channels=128.	2. Butom (orm2a(num_routares out_onamois)
	out channels=256)	3: ReLU(inplace=True)
5:	Block2D(in channels=256,	4: Conv2d(out channels, out channels,
	out_channels=512)	kernel_size=3, padding=1, bias=False)
6:	Block2D(in_channels=512,	5: BatchNorm2d(num_features=out_channels)
	out_channels=1024)	
7:	MaxPool2d(2)	6: ReLU(inplace=True)
8:	Decoder2D:	
9:	ConvTranspose2d(1024 + 5, 512, 2, stride=2)	
10:	Block2D(in_channels=1024,	
	out_channels=512)	
11:	ConvTranspose2d(512, 256, 2, stride=2)	
12:	Block2D(in_channels=512,	
	out_channels=256)	
3:	Conv Iranspose2d(256, 128, 2, stride=2)	
4:	Block2D(in_channels=256,	
15.	ConvTranspace2d(128, 64, 2, stride=2)	
15:	$\frac{120}{120} = \frac{120}{120} = $	
0.	DIOCK2D(III_CHAIIIICIS=128, Out_CHAIIIICIS=04)	
17.	ConvTranspose2d(64, 1, 2, stride-2)	
18:	Block2D(in channels=64, out channels=1)	
19:	Conv2d(64, 1, 1)	
Rol	PE NSE. We have selected the best NPE based	on the validation set with 10000 examples
gen	erated with the simulator. The NPE is fixed to or	he best-of-all model. We fine-tune the NCDE
with	a learning rate equal to 1e-4 for 2000 gradient sto	eps on on 80% of the calibration set. We use a
1-sa	mple Monte Carlo estimate of the expectation in (<u>6</u>).
I_N	PF To train I NPF we simply randomly use a bat	-
וייע. and	of 50% (θ , x) from the calibration set. We use the	he same architecture and hyper-parameters as
the	SBI NPE. The best model is selected based on the k	best training set performance. We do 50 epochs
with	1000 simulated examples for each epoch. Simulat	tions are generated randomly for each batch by
sam	pling the prior and simulating for the correspondi	ng parameters. The batch size is 100.
/		

HVAE. For the HVAE, we re-use the NPE model as the physics encoder and use the simulator as is as it is differentiable without additional effort. In addition, we follow the approach of Takeishi & Kalousis (2021) and have 1) a real-world encoder that maps \mathbf{x}_o to \mathbf{z}_a , 2) a reality-to-physics encoder, and 3) a physics-to-reality decoder. The real-world encoder has the same architecture as the NSE of the NPE and outputs the mean and log-variance of a 5D latent vector \mathbf{z}_a . The reality-to-physics and physics-to-reality also have the same architectures and are two conditional 2D U-Net with the architecture described by Algorithm 7

1396 To train the HVAE, we freeze the parameters of the NPE and optimizes the ELBO as well as a 1397 calibration loss that evaluates the likelihood assigned to the true physical parameters. All distributions 1398 are parameterized by Gaussian with mean and log-variance predicted by the neural networks. We 1399 do not use any additional losses as we expect constraining NPE and using the calibration set should 1400 already provide the necessary support to use the physics in a meaningful way. The HVAE is trained 1401 on the 2000 test examples as it is the only real-world data, calibration set aside, that we have access to. We use a batch size equal to 100 and a learning rate equal to 1e-3. We believe obtaining a better 1402 HVAE is possible. However, we emphasize the complexity of setting up a good HVAE for the only 1403 purpose of statistical inference over parameters.

1404 DATASETS

For this task, we can generate samples (θ, \mathbf{x}_s) on the fly to train the NPE. However, the calibration and test sets are real-world data. We ensure there is not overlap between calibration and test set. The is no randomization and the test set is constant for all experiments, the calibration set are also fixed for a given calibration set size.

- 1410 I.5 TASK F: WIND TUNNEL
- 1412 DESCRIPTION

1413 We use one of the wind-tunnel datasets from the causal chamber project (Gamella et al., 2024) 1414 causalchamber.org). In particular, we use the data from the load out 0.5 osr 1415 downwind_4 experiment in the wt_intake_impulse_v1 dataset. The tunnel is a chamber 1416 with two controllable fans that push air through it and barometers that measure air pressure at different 1417 locations. A hatch precisely controls the area of an additional opening to the outside (see Gamella 1418 et al. 2024, Figure 2). The data is a collection of pressure curves that result from applying a short 1419 impulse to the intake fan load and measuring the change in air pressure using one of the barometers 1420 inside the tunnel. Our inference task consists of predicting the hatch position, $\theta := [H] \in [0, 45]$ 1421 given a pressure curve (see Figure 2). As a simulator model, we combine the models A2 and C3 described in Gamella et al. (2024, Appendix D); we numerically solve the ODE in model A2, and 1422 add stochastic components to simulate the sensor noise and the unknown time point at which the 1423 impulse is applied. This results in the simulator being neither differentiable nor deterministic. A 1424 Python implementation of the complete simulator is available in the causalchamber package 1425 (models.simulator_a2_c3); visit causalchamber.org for more details. Misspecifica-1426 tion arises from the many simplifying assumptions needed to model the complex dynamics of the 1427 airflow inside the tunnel—more details can be found in Gamella et al. (2024, Appendix D.IV.1.2). 1428

Neural Posterior Estimator. The NSE and NCDE have the same 1D convolutional neural network as for Task A. For training the NPE, we use a batch size of 100 and a learning factor equal to 5e-4. NPE is trained until convergence. Other parameters are set to default values and should marginally impact the NPE obtained.

RoPE NSE. We have selected the best NPE based on the validation set with 10000 examples generated with the simulator. The NPE is fixed to one best-of-all model. We fine-tune the NCDE with a learning rate equal to 1e-4 for 20000 gradient steps on on 80% of the calibration set. We use a 1-sample Monte Carlo estimate of the expectation in (6).

J-NPE. To train J-NPE, we simply randomly use a batch composed of 50% of simulated pairs (θ, \mathbf{x}_s) and of 50% (θ, \mathbf{x}_o) from the calibration set. We use the same architecture and hyper-parameters as the SBI NPE. The best model is selected based on the best training set performance. We do 50 epochs with 10000 simulated examples for each epoch. The batch size is 100.

HVAE. There is no HVAE for this experiment as the simulator is non-differentiable.

1443 DATASETS 1444

For this task, although slightly slower than Task A and B, we can generate samples (θ, \mathbf{x}_s) on the fly to train the NPE. However, the calibration and test sets are real-world data. We ensure no overlap between the two sets for all calibration set sizes. All sets are fixed for all experiments.

- 1449
- 1450 1451
- 1452
- 1453
- 1454
- 1455
- 1456
- 1457

¹⁴⁵⁸ J COMPUTING ACAUC

Algo	prithm 9 Statistical Calibration of Posterior Distribution		
Inn	ut: Dataset of pairs $\mathcal{D} = \{(\theta^i \mathbf{x}^i)\}$ Posterior estimator $\tilde{p}(\theta \mathbf{x})$ Number of samples N		
Out	put: ACAUC		
1.	AVG CALIBRATION $= 0$		
2.	for $k \in \{1, K\}$ do		
3:	3. Initialize an empty list CredI evels		
4:	for $(\theta^i, \mathbf{x}^i) \in \mathcal{D}$ do		
5:	Initialize an empty list Samples		
6:	for $i = 1$ to M do		
7:	Sample θ^j from $\tilde{p}(\theta \mid \mathbf{x}^i)$		
8:	Append θ^j to Samples		
9:	end for		
10:	Sort Samples		
11:	Compute the rank (position in ascending order) r of θ in Samples		
12:	Set CredLevels = $\frac{r}{N}$		
13:	Append CredLevel to CredLevels		
14:	end for		
15:	Sort CredLevels		
16:	CALIBRATION = $\sum_{i=1}^{N}$ CredLevels $[i] - \frac{i}{N}$		
17:	$AVG_CALIBRATION = AVG_CALIBRATION + \frac{CALIBRATION}{K}$		
18:	end for		
Ret	urn: AVG CALIBRATION		

1483 K ADDITIONAL RESULTS1484





Figure 7: Three corner plots for task A with a calibration set with 50 samples.





Figure 12: Three corner plots for task E with distribution shift with a calibration set with 50 samples.



Figure 13: Calibration plots of the different methods on the 6 benchmarks, the coverage at each level is the average of the coverage of the marginal distributions. Each color indicates a different algorithm and the opacity is proportional to the size of the calibration set which ranges from 10 to 1000. We observe that RoPE and OT-only are consistently well calibrated for.