

# Fuse It or Lose It: Deep Fusion for Multimodal Simulation-Based Inference

Marvin Schmitt\*, Leona Odole†, Stefan T. Radev‡ and Paul-Christian Bürkner†

\*University of Stuttgart, †TU Dortmund University, ‡Rensselaer Polytechnic Institute,

**Abstract**—We present multimodal neural posterior estimation (MultiNPE), a method to integrate heterogeneous data from different sources in simulation-based inference with neural networks. Inspired by advances in deep fusion, it allows researchers to analyze data from different domains and infer the parameters of complex mathematical models with increased accuracy. We consider three fusion approaches for MultiNPE (early, late, hybrid) and evaluate their performance in three challenging experiments. MultiNPE not only outperforms single-source baselines on a reference task, but also achieves superior inference on scientific models from cognitive neuroscience and cardiology. We systematically investigate the impact of partially missing data on the different fusion strategies. Across our experiments, late and hybrid fusion techniques emerge as the methods of choice for practical applications of multimodal simulation-based inference.

**Index Terms**—Machine Learning for Science, Inverse Problems, Neural Density Estimation, Uncertainty Quantification.

## I. INTRODUCTION

**S**IMULATIONS have become a fundamental tool to model complex phenomena across the sciences and engineering [1]. For example, in precision medicine, high-fidelity hemodynamics simulators mimic the blood flow through the human body. In such a simulation model, latent *parameters*  $\theta$  determine the behavior of the complex system, which outputs observable data  $\mathcal{D}$ . The latent parameters of the hemodynamics simulator, for example, are the cardiovascular characteristics of the patient, such as the left-ventricular ejection time (LVET) or the arterial diameter (see **Experiment 3** for details). The observable data in this example are a patient’s pulse waves which can easily be measured at the patient’s fingertip or wrist with medical measurement devices. Crucially, the latent parameters  $\theta$  are not directly observable without intrusive means, but they are relevant for medical practitioners who want to evaluate the patient’s cardiovascular health. Thus, we seek to *infer* the latent parameters  $\theta$  based on the observable data  $\mathcal{D}$ . The probabilistic (Bayesian) approach to this *inverse problem* leads to the posterior distribution  $p(\theta | \mathcal{D}) \propto p(\theta) p(\mathcal{D} | \theta)$ , which describes the distribution of plausible parameter values  $\theta$  given a prior belief  $p(\theta)$  and observable data  $\mathcal{D}$ .

There exists a myriad of methods to approximate the posterior distribution in the methodological repertoire of Bayesian statistics [2, 3]. However, there is a critical feature that complicates Bayesian inference on simulators. By design, it is typically straightforward to generate synthetic data from a simulation model. Yet, the observation model  $p(\mathcal{D} | \theta)$  necessary to compute the posterior distribution might be only *implicitly* defined, lacking a closed-form likelihood function [4, 5].

Implicit models cannot easily be estimated with established Bayesian algorithms like Markov chain Monte Carlo (MCMC; [6]) or variational inference (VI; [7]) to approximate the posterior distribution. Furthermore, MCMC or VI algorithms need to be re-run from scratch for every new observed data set, which makes real-time estimation or monitoring impossible.

Fueled by recent advances in generative neural networks, *amortized simulation-based Bayesian inference* solves both problems simultaneously because it (i) does not require explicit likelihoods; and (ii) yields near-instant approximate posterior draws for any new data set. More concretely, the family of neural posterior estimation (NPE) algorithms directly learns a surrogate posterior  $q_\phi(\theta | \mathcal{D}) \approx p(\theta | \mathcal{D})$  from simulations of the joint model  $p(\theta, \mathcal{D})$  via neural network training (see **Section II**). Subsequently, the upfront training is *amortized* by rapid posterior inference: For a *new* observed data set  $\mathcal{D}_o$ , the neural network can instantly generate draws from the approximate posterior  $q_\phi(\theta | \mathcal{D}_o)$ , making real-time Bayesian inference feasible for a large class of applied problems.

Amortized simulation-based inference is still in its infancy, and we extend its repertoire to the practically relevant class of mechanistic *multimodal models*, where a set of shared parameters influences heterogeneous data sources via distinct simulators. Returning to the running example of computational models in cardiovascular precision medicine, we often use different measurement devices to monitor the pulse waves of a single patient [8], and doctors have additional behavioral or demographical data on patients. Naturally, we want to *combine* all this information into a holistic analysis that accounts for all available data on the patient’s cardiovascular health. Neural simulation-based inference currently lacks the tools to properly analyze such multimodal data. Our paper addresses this limitation with the following contributions (see **Figure 1**):

- 1) We present multimodal neural posterior estimation (MultiNPE), which enables the integration of multimodal data into amortized simulation-based inference methods.
- 2) We develop variations of MultiNPE, translating advances in attention-based deep fusion learning into probabilistic machine learning with neural networks.
- 3) We demonstrate that MultiNPE outperforms existing simulation-based inference methods on a 10-dimensional reference task as well as two applied scientific problems from cognitive neuroscience and cardiology.

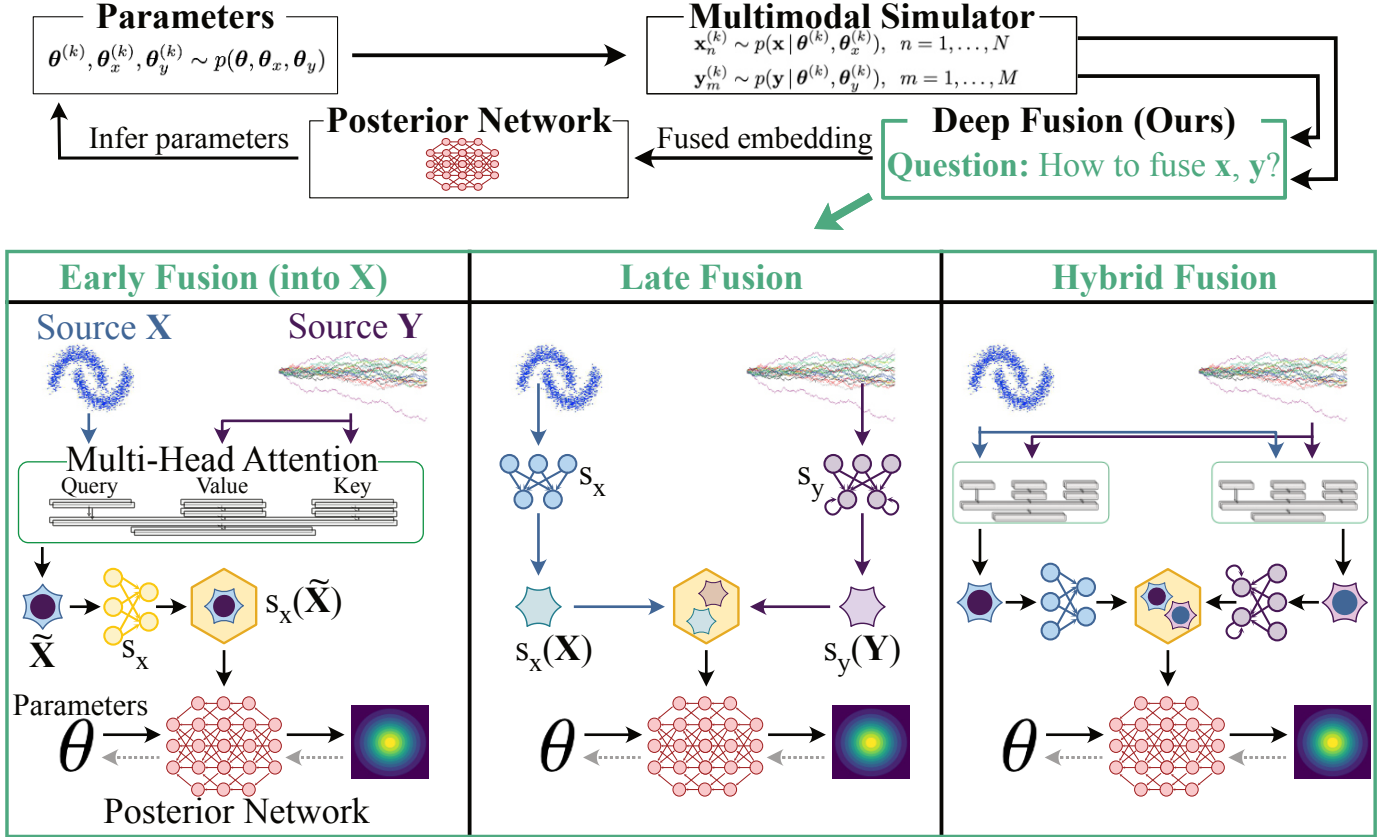


Fig. 1: We present a set of deep fusion methods to equip simulation-based inference (SBI) with the ability to integrate information from multiple heterogeneous data sources. Early fusion (into  $\mathbf{X}$ ) uses multi-head attention with  $\mathbf{X}$  as query and  $\mathbf{Y}$  as both key and value, yielding the cross-informed representation  $\tilde{\mathbf{X}}$ , followed by another summary network  $s_x$ . In contrast, late fusion learns separate embeddings  $s_x(\mathbf{X})$  and  $s_y(\mathbf{Y})$ , and then fuses the embeddings. Hybrid fusion combines both worlds by using cross-shaped multi-head attention like in early fusion, followed by separate embeddings and a late fusion step. See [Section III](#) for a more formal specification.

## II. PRELIMINARIES

This section gives a brief overview of neural posterior estimation, learned summary statistics, and multi-head attention. Acquainted readers can fast-forward to [Section III](#).

### A. Neural posterior estimation

The inverse problem of approximating the posterior distribution in simulation-based inference (SBI) can be tackled directly with a class of algorithms called neural posterior estimation (NPE). NPE uses a neural network to estimate the surrogate conditional density  $q_\phi(\theta | \mathcal{D})$ , where  $\phi$  denotes learnable neural network weights. Common neural network architectures include normalizing flows [9], score-based diffusion [10], flow matching [11], or consistency models [12]. Here, we mainly focus on normalizing flows due to their fast single-pass inference and simple training, even though our approach for multimodal SBI translates seamlessly to other backbone conditional neural density estimators (see [Experiment 3](#) for a demonstration of flow matching).

A conditional normalizing flow learns a bijective map between a simple base distribution (e.g., a unit Gaussian) and the target posterior  $p(\theta | \mathcal{D})$ . The normalizing flow is

optimized by minimizing the Kullback-Leibler (KL) divergence between the true posterior  $p(\theta | \mathcal{D})$  and its approximation  $q_\phi(\theta | \mathcal{D})$  via the maximum likelihood objective  $\mathbb{E}_{p(\theta) p(\mathcal{D} | \theta)} [-\log q_\phi(\theta | \mathcal{D})]$ .<sup>1</sup> The training data for the normalizing flow are synthetic tuples  $(\theta, \mathcal{D})$  which are generated through an ancestral sampling procedure

$$\begin{aligned} \theta &\sim p(\theta) \\ \mathcal{D} &\sim p(\mathcal{D} | \theta), \end{aligned} \quad (1)$$

arising from factorizing the joint distribution  $p(\theta, \mathcal{D})$  into the prior  $p(\theta)$  and the (multimodal) observation model  $p(\mathcal{D} | \theta)$ . Once the normalizing flow has been trained with a maximum likelihood objective (see [Equation 2](#) below), it can instantly sample from the posterior  $q_\phi(\theta | \mathcal{D}_o)$  for new observed data  $\mathcal{D}_o$ . Thus, by re-casting costly probabilistic inference as a neural network prediction task, normalizing flows achieve

<sup>1</sup>While the optimization objectives of variational inference and normalizing flows look strikingly similar, they differ in a fundamental aspect: Variational inference optimizes the *reverse* KL divergence, which in turn requires access to the joint density of the model and leads to mode-seeking behavior. In contrast, normalizing flows target the *forward* KL divergence. As a consequence, they do not require access to the joint density and can be trained in a fully simulation-based (aka. likelihood-free) setting. Further, this generally leads to mode-covering behavior.

*amortized inference* across the space of typical samples from the joint model  $p(\theta, \mathcal{D})$ .

### B. Embedding networks for end-to-end learned summary statistics

In Bayesian inference, the data  $\mathcal{D}$  can be replaced by *sufficient* summary statistics  $s_*(\mathcal{D})$  without altering the posterior:  $p(\theta | \mathcal{D}) = p(\theta | s_*(\mathcal{D}))$ . Ideally,  $s_*$  is also low-dimensional, achieving lossless compression with respect to  $\theta$  conditioned on  $\mathcal{D}$ . While low-dimensional sufficient summary statistics are notoriously difficult to find for complex problems, the task of constructing approximate summary statistics  $s(\mathcal{D})$  with  $p(\theta | \mathcal{D}) \approx p(\theta | s(\mathcal{D}))$  has been extensively studied for non-neural approximate Bayesian inference [13, 14, 15, 16]. Within neural SBI, specialized neural networks are employed to learn embeddings of the data  $\mathcal{D}$  in tandem with the posterior approximator [17, 18, 19, 20]. These embedding networks  $s_\psi$  learn a transformation that aims to obtain low-dimensional statistics of the data  $\mathcal{D}$  which are sufficient for posterior inference (not necessarily for reconstructing the data). The embedding networks are parameterized by learnable neural network weights  $\psi$ . The NPE loss with *learned embeddings*  $s_\psi(\mathcal{D})$  minimizes the maximum likelihood objective

$$\mathcal{L}(\phi, \psi) = \mathbb{E}_{p(\theta)p(\mathcal{D}|\theta)} [-\log q_\phi(\theta | s_\psi(\mathcal{D}))], \quad (2)$$

and we omit the network weights  $\psi$  for brevity in the following. The concrete architecture of the embedding network should match the probabilistic symmetries of the data. For example, *i.i.d.* data sets can be embedded with a permutation-invariant neural network, such as a DeepSet [21] or a Set-Transformer [22]. Similarly, time series data require a neural architecture which respects their temporal dependencies, such as an LSTM [23] or a temporal fusion transformer [24].

### C. Multi-head attention

Attention mechanisms play a crucial role in deep learning, and one of the most notable architectures that has taken the field by a storm is the Transformer [25]. The Transformer introduces a highly effective mechanism for capturing dependencies and relationships within sequences of data, making it particularly well-suited for tasks such as natural language processing [25] or computer vision [26]. The core of the Transformer’s attention mechanism is the scaled dot-product attention, defined as

$$\text{Attention}(Q, V, K) = \text{softmax}\left(\frac{QK^\top}{\sqrt{d_K}}\right)V, \quad (3)$$

with queries  $Q$ , values  $V$ , and keys  $K$  of dimension  $d_K$ . To enhance the model’s ability to capture different types of relationships and dependencies in the data, the Transformer employs multi-head attention (MHA). MHA enables the model to jointly attend to information from different subspaces of the data across multiple attention heads. Each attention head is a separate instance of the scaled dot-product attention mechanism (Equation 3), and their outputs are combined using

learnable linear transformations to produce the final multi-head attention output,

$$\begin{aligned} \text{MHA}(Q, V, K) &= [\text{head}_1, \dots, \text{head}_h] W^O \\ \text{head}_i &= \text{Attention}(QW_i^Q, KW_i^K, VW_i^V) \end{aligned} \quad (4)$$

where  $h$  represents the number of attention heads, and  $W^O$ ,  $W_i^Q$ ,  $W_i^K$ ,  $W_i^V$  are learnable weight matrices [25]. The multi-head attention mechanism allows the Transformer model to encode patterns and relationships in the data via the matrix  $W^O$ , which makes the architecture highly effective for a range of sequence-to-sequence tasks.

## III. METHOD

### A. Simulation paradigm and notation

In this section, we consider multimodal<sup>2</sup> test data  $\mathcal{D}_o = \{\mathbf{X}_o, \mathbf{Y}_o\}$  from two sources<sup>3</sup>, as well as a simulation program capable of generating synthetic data  $\mathcal{D} = \{\mathbf{X}, \mathbf{Y}\}$ . An instance of either data source can consist of multiple observations (e.g., patients in a clinical trial) or discrete steps in a time series. We use  $N$  for the cardinality of the first source,  $\mathbf{X} \equiv \{\mathbf{x}_n\}_{n=1}^N$ , and  $M$  for the cardinality of the second source  $\mathbf{Y} \equiv \{\mathbf{y}_m\}_{m=1}^M$ . Following the standard notation in SBI, the neural networks are trained on a total of  $K$  data sets  $\{\mathcal{D}^{(k)}\}_{k=1}^K \equiv \{\{\mathbf{X}^{(k)}, \mathbf{Y}^{(k)}\}\}_{k=1}^K$ .  $K$  is also called the *simulation budget*. To shorten notation, we drop the data set index  $k$  when it is clear from the context. The sub-programs generating the individual data modalities are based on common parameters  $\theta$  as well as domain-specific parameters  $\theta_x, \theta_y$ . Using the verbose notation once to avoid ambiguity, the joint forward model  $p(\theta, \theta_x, \theta_y, \mathbf{X}, \mathbf{Y})$  for a single data set  $\mathcal{D}^{(k)} = \{\mathbf{X}^{(k)}, \mathbf{Y}^{(k)}\}$  is defined as:

$$\begin{aligned} \theta^{(k)}, \theta_x^{(k)}, \theta_y^{(k)} &\sim p(\theta, \theta_x, \theta_y) \\ \mathbf{x}_n^{(k)} &\sim p(\mathbf{x} | \theta^{(k)}, \theta_x^{(k)}), \quad n = 1, \dots, N \\ \mathbf{y}_m^{(k)} &\sim p(\mathbf{y} | \theta^{(k)}, \theta_y^{(k)}), \quad m = 1, \dots, M \end{aligned} \quad (5)$$

The result of sampling from this forward model  $K$  times is a set of  $K$  tuples of parameters and data sets,

$$\left\{ \underbrace{\{\theta^{(k)}, \theta_x^{(k)}, \theta_y^{(k)}\}}_{\text{Parameters}}, \underbrace{\{\mathbf{X}^{(k)}, \mathbf{Y}^{(k)}\}}_{\text{Data } \mathcal{D}^{(k)}} \right\}_{k=1}^K,$$

and the inverse problem consists of inferring all unknown parameters from the data. Since SBI relies on synthetic data, the ground-truth parameter values are known and available during the training phase. In the inference (test) phase however, the ground truth parameters of the test data  $\mathcal{D}_o = \{\mathbf{X}_o, \mathbf{Y}_o\}$  are unknown and need to be estimated by the generative network  $q_\phi(\theta, \theta_x, \theta_y | \mathbf{X}_o, \mathbf{Y}_o)$ .

<sup>2</sup>Within the scope of this paper, we use the term *multimodal* generally for any two datasets whose structure would make them incompatible for neural posterior estimation because they are non-concatenable and common tricks (e.g., zero-padding) are invalid. For example, in **Experiment 1** we combine Gaussian *i.i.d.* data  $\mathbf{X}$  and time series data  $\mathbf{Y}$ .

<sup>3</sup>We limit this description to two sources for brevity. As discussed in **Section III-D** and illustrated in **Experiment 3**, a more involved layout of attention blocks can readily fuse more sources in a similar fashion.

### B. Necessity of principled deep fusion

Consider the scenario where we estimate a single shared parameter  $\theta$  that manifests itself in both an *i.i.d.* data set  $\mathbf{X} \sim p(\mathbf{X}|\theta)$  and a Markovian time series  $\mathbf{Y} \sim p(\mathbf{Y}|\theta)$ . The fundamentally different probabilistic systems for  $\mathbf{X}$  and  $\mathbf{Y}$  cannot possibly be efficiently learned with a single neural architecture because (i) a permutation invariant network is suited for *i.i.d.* data but cannot capture the autoregressive structure of a time series; and (ii) a time series network can fit time series but cannot efficiently learn the permutation-invariant structure of *i.i.d.* data. As a consequence, we need separate information processing streams to accommodate the specific structure of each data source. Yet, the neural density estimator  $q_\phi$  demands a fixed-length conditioning vector  $s(\mathcal{D})$ . We serve both requirements simultaneously: First, we process the heterogeneous streams of information  $\mathbf{X}$  and  $\mathbf{Y}$  with dedicated architectures. Second, we integrate the processed information into a fixed-length embedding before it enters the neural density estimator.

When either of the simulators has no individual parameters, we arrive at a special case of Equation 5, where one simulator only features shared parameters. This has no effect on our method, which remains unaltered (see **Experiments 1** and **3**). However, if there are no *shared* parameters at all, a multimodal architecture will clearly have no advantage over separate inference algorithms for the two sub-problems because there is no mutual information to leverage via weight sharing.

### C. Fusion strategies

The integration of information from different data sources is called *fusion*, and there are multitudes of options for *how* and *when* the fusion happens (see Figure 1). Previous work on deep fusion learning differentiates early fusion, late fusion, and hybrid approaches [27, 28, 29]. Our embedding network  $s(\cdot)$  corresponds to the *decision level function* in the standard multimodal machine learning taxonomy.

**Early fusion** performs the fusion step as early as possible, ideally directly on the raw data (see Figure 1, panel 1). We implement this via cross-attention [30, 31] between the input modalities  $\mathbf{X}$  and  $\mathbf{Y}$ . Concretely, we use multi-head attention [25] with queries  $Q$ , values  $V$ , and keys  $K$ . In multi-head attention, the data inputs  $\mathbf{X}$  and  $\mathbf{Y}$  can differ with respect to their dimensions but the shapes of  $V$  and  $K$  must align. Thus, we select one of the data sources as query  $Q$  while the other one acts as both value  $V$  and key  $K$ . In **Experiment 1**, we illustrate that the choice of data sources for  $Q$  and  $V, K$  is important. After the attention-based fusion step, we pass the output of the multi-head attention block to an appropriate embedding network  $s(\cdot)$  to provide a fixed-length input for the conditional neural density estimator  $q_\phi$ . In summary, the information flow in early fusion is formalized as

$$\begin{aligned} \text{Early Fusion to } \mathbf{X}: s(\mathcal{D}) &= s_x\left(\text{MHA}(Q(\mathbf{X}), V(\mathbf{Y}), K(\mathbf{Y}))\right), \\ \text{Early Fusion to } \mathbf{Y}: s(\mathcal{D}) &= s_y\left(\text{MHA}(Q(\mathbf{Y}), V(\mathbf{X}), K(\mathbf{X}))\right), \end{aligned} \quad (6)$$

where  $\text{MHA}(Q, V, K)$  denotes multi-head attention (see Section II-C for details).

**Late fusion** introduces the fusion step at a later stage (see Figure 1 panel 3). In SBI with learnable embeddings, this translates to fusion immediately before passing the final embedding to the conditional neural density estimator as conditioning variables. At this stage, both data inputs have already been embedded into learned summary statistics  $s_x(\mathbf{X})$  and  $s_y(\mathbf{Y})$ . Thus, late fusion can be achieved by simply concatenating the embeddings,  $s(\mathcal{D}) = g(s_x(\mathbf{X}), s_y(\mathbf{Y})) = [s_x(\mathbf{X}), s_y(\mathbf{Y})]$ , which is a common choice for the fusion function  $g$  [27, 28, 29].

**Hybrid fusion** combines early and late fusion (see Figure 1 panel 4). Initially, we use cross attention with *both*  $\mathbf{X}$  and  $\mathbf{Y}$  as the query  $Q$ : We construct a cross-shaped information flow where we embed each data source using cross-attention information from the other source. This leads to a *symmetrical* information flow and overcomes the drawback of early fusion, where one of the data sources must be chosen as the query  $Q$ . The outputs of the symmetrical cross-attention step are then each passed to an embedding network  $s_x(\cdot), s_y(\cdot)$ , and the information streams are fused just before entering the neural density estimator:

$$\begin{aligned} \tilde{\mathbf{X}} &= \text{MHA}(Q(\mathbf{X}), V(\mathbf{Y}), K(\mathbf{Y})) \\ \tilde{\mathbf{Y}} &= \text{MHA}(Q(\mathbf{Y}), V(\mathbf{X}), K(\mathbf{X})) \\ s(\mathcal{D}) &= g(s_x(\tilde{\mathbf{X}}), s_y(\tilde{\mathbf{Y}})) = [s_x(\tilde{\mathbf{X}}), s_y(\tilde{\mathbf{Y}})] \end{aligned} \quad (7)$$

We hypothesize that hybrid fusion enables more flexible resource allocation: Features of an informative source as well as interactions can be captured in the embedding network, which reduces the burden on the generative network  $q_\phi$ .

### D. More than two data sources

This section will discuss the natural extension of our fusion architectures beyond two sources. In the following, let  $L \in \mathbb{N}_{\geq 2}$  be the number of data sources  $\mathcal{D} = \{\mathcal{D}_l\}_{l=1}^L$ .

Late fusion naturally translates to an arbitrary number of sources: Each source  $\mathcal{D}_l$  has a dedicated embedding network  $s_l(\mathcal{D}_l)$  to learn sufficient summary statistics for posterior inference. Finally, all embeddings are combined into a joint embedding  $s(\mathcal{D}) = g(s_1(\mathcal{D}_1), \dots, s_L(\mathcal{D}_L))$  with suitable  $g$  (e.g., concatenation as above). Thus, the number of networks in late fusion scales linearly in  $\mathcal{O}(L)$ . Early and hybrid fusion, however, involve pairwise cross-attention blocks, which do not natively generalize to  $L \geq 2$  inputs.

For early fusion, there are  $L!$  options to choose the layout of pairwise cross-attention fusion blocks, but only  $1 + \dots + (L-1)$  blocks must be realized in practice to implement a cascade of cross-attention steps for early fusion. In addition, we require one embedding network, leading to a total of  $1 + 2 + \dots + (L-1) + 1 \in \mathcal{O}(L^2)$  networks.

In hybrid fusion, however, we want a full cross-exchange of information across all sources, which requires a total of  $L!$  networks. In addition, each source needs one embedding network. This leads to a total of  $L! + L \in \mathcal{O}(L!)$  networks, which clearly raises scaling issues for large  $L$ . In all three experiments, we carefully compare whether the less scalable hybrid fusion approach yields a worthwhile advantage over the



more scalable late fusion method, and we will conclude that late fusion is a viable option in most applications.

#### IV. RELATED WORK

**Multimodal fusion.** Researchers have long been integrating different types of features to improve the performance of machine learning systems [32]. As [33] remark, using deep fusion to learn *fused representations of heterogeneous features* in multimodal settings is a natural extension of this strategy. Recently, transformers have been employed for multimodal problems across many applications [34], such as image and sentence matching [35], multispectral object detection [36], or integration of image and depth information [37]. We confirm the potential of cross-attention in probabilistic machine learning with conditional neural density estimators. All of our fusion variants implement unified embeddings for heterogeneous data sources, corresponding to *joint representation* in the taxonomy of [33].

**Multimodality and missing data.** Multimodal learning algorithms can naturally address the problem of missing data because missing information from one source may be compensated for by another source (see **Experiment 2**). In the context of multimodal time series, this has been addressed via factorized inference on state space models [38] and learned representations via tensor rank regularization [39]. Our multimodal NPE method also learns robust representations from partially missing data, but we use fusion techniques that respect the probabilistic symmetry of the data, rather than a certain factorization of the posterior distribution. Bayesian meta-learning [40] has been used to study the efficiency of multimodal learning under missing data, both during training and inference time [41]. Similarly, our approach embodies Bayesian meta-learning principles by extending the *amortization scope* of NPE to incorporate missing data [42], which is in turn facilitated by our fusion schemes.

**Hierarchical Bayesian models.** Hierarchical or multilevel Bayesian models [43, 44] are used to model the dependencies in nested data, where observations are organized into clusters or levels. While these models often feature *shared* parameters across observational units or *global* parameters describing between-cluster variations [45], they focus on analyzing the variations of a *single data modality at different levels*. In contrast, multimodal models capitalize on integrating information from different sources or modalities. That being said, a multimodal problem could also be formulated in a hierarchical way, such that the shared parameters of different modalities admit a hierarchical prior. While the complexity of such models quickly becomes prohibitive, our MultiNPE approach could pave the way for hierarchical multimodal approaches where the latter have been foregone merely out of computational desperation.

#### V. EMPIRICAL EVALUATION

**Settings.** We evaluate MultiNPE in a synthetic multimodal model with fully overlapping parameter spaces across the data modalities (**Experiment 1**), a neurocognitive model with partially overlapping parameter spaces and missing data (**Experiment 2**), and a cardiovascular data set with three data sources (**Experiment 3**).

**Evaluation metrics.** For all experiments, we evaluate the accuracy of the posterior estimates as well as their uncertainty calibration and Bayesian information gain on  $J$  unseen test data sets  $\{\mathcal{D}_o^{(j)}\}_{j=1}^J$  with known ground-truth parameters  $\{\theta_*^{(j)}\}_{j=1}^J$ . In **Experiment 1**, we additionally report the distance between our approximate posterior and a reference ground-truth posterior via the maximum mean discrepancy (MMD).<sup>4</sup> Let  $\{\theta_s^{(j)}\}_{s=1}^S$  be the set of  $S$  posterior draws from the neural approximator  $q_\phi(\theta | \mathcal{D}_o^{(j)})$  conditioned on the data set  $\mathcal{D}_o^{(j)}$ . To quantify accuracy, we compute the average root mean square error (RMSE) between posterior draws and ground truth parameter values over the test set:

$$\text{RMSE} = \frac{1}{J} \sum_{j=1}^J \sqrt{\frac{1}{S} \sum_{s=1}^S (\theta_s^{(j)} - \theta_*^{(j)})^2} \quad (8)$$

We quantify uncertainty calibration via simulation-based calibration (SBC; [47]): in proper Bayesian inference, all regions  $U_q(\theta | \mathcal{D})$  of the *true* posterior  $p(\theta | \mathcal{D})$  are well calibrated for any quantile  $q \in (0, 1)$  by design [48], that is, the equality

$$q = \iint \mathbf{I}[\theta_* \in U_q(\theta | \mathcal{D})] p(\mathcal{D} | \theta_*) p(\theta_*) d\theta_* d\mathcal{D} \quad (9)$$

always holds, where  $\mathbf{I}[\cdot]$  is the indicator function. Discrepancies from this equality indicate deficient calibration of the approximate posterior. We report the median SBC error of central credible intervals computed for 20 linearly spaced quantiles  $q \in [0.5\%, 99.5\%]$ , averaged across the test set (i.e., expected calibration error; ECE). Third, we quantify the (Bayesian) information gain via the posterior contraction based on the average ratio between posterior and prior variance across the test data,

$$\text{Contraction} = \frac{1}{J} \sum_{j=1}^J \left( 1 - \frac{\text{Var}_\theta [p(\theta | \mathcal{D}^{(j)})]}{\text{Var}_\theta [p(\theta)]} \right). \quad (10)$$

Finally, we use the maximum mean discrepancy (MMD) to quantify the distance between the approximate and ground-truth posterior distribution based on samples [49]. All four metrics are *global*: they estimate performance across the entire joint model  $p(\theta, \mathcal{D})$  instead of singling out particular data sets or true model parameters [48]. The metrics can directly be computed based on test simulations from the joint model, which is essentially instantaneous due to amortized inference.

<sup>4</sup>This is possible because **Experiment 1** entails a likelihood-based model, allowing for posterior sampling with gold-standard Hamiltonian Monte Carlo, as implemented in the Stan software [46].

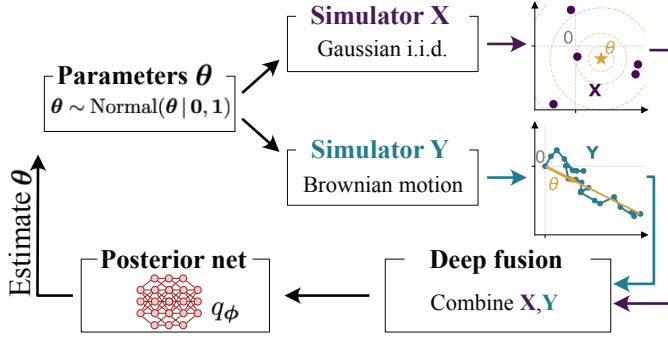


Fig. 2: **Experiment 1:** Simplified 2D visualization of the experimental setup. The actual experiment is implemented in 10-dimensional spaces for both the parameters  $\theta$  and the observed measurement variables  $\mathbf{X}, \mathbf{Y}$ .

#### A. Experiment 1: Exchangeable data and Brownian motion

This experiment compares MultiNPE and standard NPE on a synthetic task where a common parameter  $\theta \in \mathbb{R}^{10}$  is used as (i) the location parameter of Gaussian *i.i.d.* data  $\mathbf{X} \in \mathbb{R}^{5 \times 10}$  and (ii) the drift rate of a stochastic trajectory  $\mathbf{Y} \in \mathbb{R}^{20 \times 10}$ ,

$$\begin{aligned} \theta &\sim \text{Normal}(\theta | \mathbf{0}, \mathbf{1}), \\ \mathbf{x}_n &\sim \text{Normal}(\mathbf{x} | \theta, \mathbf{1}), \quad n = 1, \dots, N \\ d\mathbf{y}_m(t) &= \theta dt + \sigma d\mathbf{W}(t), \quad m = 1, \dots, M \\ &\text{with } \mathbf{W}(t) \sim \text{Normal}(\mathbf{W} | \mathbf{0}, \mathbf{1}), \end{aligned} \quad (11)$$

where the *i.i.d.* data consist of  $N = 5$  observations  $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_5\}$ , the trajectory is discretized into  $M = 20$  steps  $\mathbf{Y} = (\mathbf{y}_1, \dots, \mathbf{y}_{20})$  with noise  $\sigma = 0.5$ , an interval of  $t = [0, 3]$ , and the initial condition is  $\mathbf{y}_1 = \mathbf{0}$ .

We compare the following neural approximators: NPE with input  $\mathbf{X}$ , NPE with input  $\mathbf{Y}$ , as well as MultiNPE variants with early fusion to  $\mathbf{X}$ , early fusion to  $\mathbf{Y}$ , late fusion, and hybrid fusion. Each neural approximator is trained on the same training set with a simulation budget of  $K = 5000$  for 30 epochs, and we repeat each training and evaluation process ten times. All data originating from the *i.i.d.* source ( $\mathbf{X}$  or  $\tilde{\mathbf{X}}$ ) are embedded with a set transformer [22], and data on the time series stream ( $\mathbf{Y}$  or  $\tilde{\mathbf{Y}}$ ) are embedded with a temporal fusion transformer [24].

**Results.** We repeat each neural network training ten times with different random number generator seed and base our evaluations on 1000 unseen test data sets, for each of which we draw 1000 posterior samples per architecture (6 architectures) and training repetition (10 repetitions). As a consequence, our systematic evaluation is based on a total of 60 million approximate posterior samples. We observe that late fusion and hybrid fusion outperform standard NPE architectures which only have access to a single data source (see Figure 3 and Table I), as evidenced by lower RMSE, lower expected calibration error (ECE), higher posterior contraction, and lower MMD to a reference posterior. It is evident that data source  $\mathbf{X}$  is less informative for posterior inference than data source  $\mathbf{Y}$  (“only  $\mathbf{X}$ ” performs much worse than “only  $\mathbf{Y}$ ”). As a consequence, early fusion to  $\mathbf{X}$  leads to worse performance than early fusion to  $\mathbf{Y}$ . We conclude that the efficacy of early

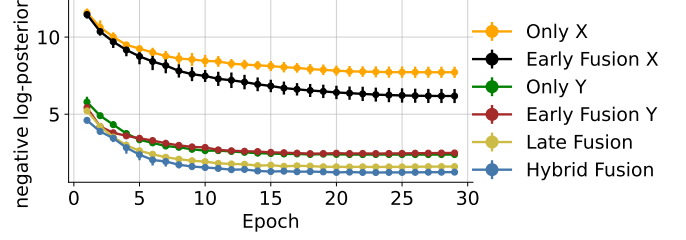


Fig. 3: **Experiment 1:** Two of our multimodal schemes (late fusion and hybrid fusion) outperform single-source architectures (only  $\mathbf{X}/\mathbf{Y}$ ), as indexed by better (lower) negative log posterior on held-out data across ten repetitions with different seeds.

fusion depends on the informativeness of the data source used as a query. Knowing which source contains more information a priori is a hyperparameter design choice that we would ideally like to avoid in real-world applications. In fact, neither late nor hybrid fusion require such a design choice, and both schemes outperform early fusion in this experiment at the cost of longer neural network training.<sup>5</sup> The most expressive neural architecture, hybrid fusion, shows the best performance by a small margin. Heuristically, it combines the best of both worlds: Hybrid fusion extracts information from the raw data by the  $\mathbf{X}$ -shaped cross-attention modules but avoids the necessity of choosing one of the domains to early-fuse into. Yet, the performance gain over late fusion is small; thus, late fusion might be employed in scenarios where practical considerations (e.g., many sources, limited time) prohibit a hybrid approach.

<sup>5</sup>While the training time scales linearly with the number of multi-head attention modules, the summary networks are also based on self-attention. Thus, *runtime complexity* is unaffected by our additional fusion scheme.

Architecture	Time <sup>1</sup> ↓	RMSE ↓	ECE [%] ↓	Contraction ↑	MMD ↓
Only $\mathbf{X}$	117 (110, 149)	0.81 (0.80, 0.89)	1.43 (0.98, 1.84)	0.68 (0.61, 0.68)	1.89 (0.035)
Only $\mathbf{Y}$	<b>100</b> (95, 141)	0.40 (0.40, 0.40)	3.44 (3.02, 3.63)	0.93 (0.93, 0.93)	0.40 (0.001)
Early Fusion $\mathbf{X}$	140 (131, 150)	0.88 (0.82, 0.93)	<b>1.35</b> (1.04, 1.80)	0.61 (0.57, 0.66)	2.03 (0.050)
Early Fusion $\mathbf{Y}$	128 (118, 152)	0.45 (0.45, 0.45)	5.45 (5.06, 5.91)	0.91 (0.91, 0.91)	0.63 (0.002)
Late Fusion	193 (172, 235)	0.36 (0.36, 0.36)	4.73 (4.31, 5.21)	0.94 (0.94, 0.94)	0.28 (0.003)
Hybrid Fusion	227 (211, 280)	<b>0.35</b> (0.35, 0.35)	4.99 (4.44, 5.18)	<b>0.95</b> (0.95, 0.95)	<b>0.23</b> (0.002)

TABLE I: **Experiment 1:** Our multimodal NPE architectures are superior to single-source NPE algorithms on 1 000 unseen data sets, as indexed by improved accuracy (RMSE), information gain (contraction), and similarity to a reference posterior (MMD). The subpar calibration under  $\mathbf{Y}$  propagates into the fused posteriors. The table shows median (min, max) across ten training runs of each architecture for time, RMSE, ECE, and contraction; and mean (SE) across training runs for MMD.

<sup>1</sup> Training time [seconds]

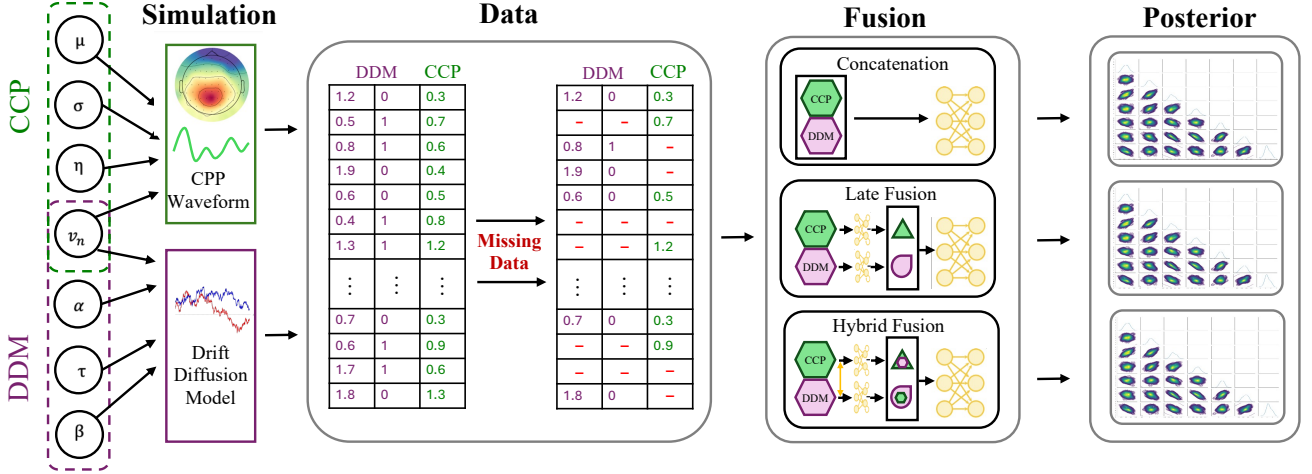


Fig. 4: **Experiment 2.** Overview of the experimental setup. A human’s neurocognitive attributes parameterize the simulation programs for **centro-parietal positivity (CPP)** and **reaction times (DDM)**.

### B. Experiment 2: Neurocognitive model of decision making and EEG with missing data

In an attempt to better understand human cognition, researchers use increasingly complex models to understand the neurological relationship between cognitive and physical phenomena. This research is a promising avenue to gain insights into the fundamental mechanisms underlying human information processing. A human’s decision and reaction time can be modeled as a stochastic evidence accumulation process via a drift-diffusion model (DDM; for a detailed description, see [50]). In a nutshell, a DDM models human decision making as a random walk with a drift (i.e., global direction) that corresponds to the person’s information processing speed. Further, the DDM estimates (i) the information that is required to form a decision; (ii) cognitive biases that shift the starting point; and (iii) a non-decision time that accounts for purely motorical latencies (e.g., moving the hand to press a button). In addition to the DDM model, the centro-parietal positive (CPP) waveform is a neurophysiological marker associated with human decision making [51].

This experiment uses a multimodal neurocognitive model to integrate both the cognitive drift-diffusion model for decision making and an observable representation of the CPP waveforms on an EEG (Model 7 from [51]). As argued in detail by Ghaderi-Kangavari et al. [51], models that jointly integrate neural processes on a trial-level represent the state-of-the-art in cognitive modeling research to holistically represent human behavior. In this experiment, we apply our MultiNPE method to improve the probabilistic estimation in single-trial joint integrative models, which contributes to a line of research towards scalable probabilistic modeling of human behavior.

The joint cognitive forward model is characterized by six parameters (defined below) that govern two partially entangled data generating processes on a trial<sup>6</sup> level, with shared *global*

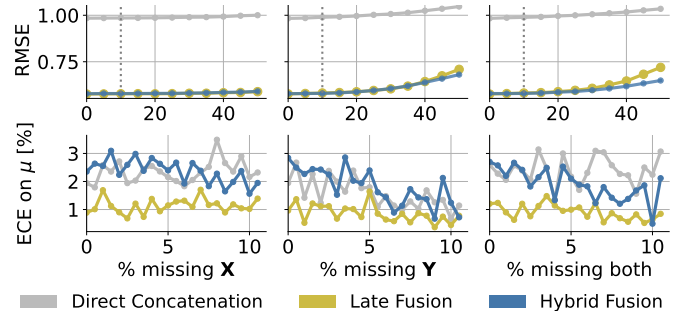


Fig. 5: **Experiment 2:** Hybrid fusion and late fusion consistently show better accuracy (RMSE averaged over all parameters) than the default (direct concatenation). Recall that the training uses 10% missing data (top row, dotted line) and missingness beyond 10% is a substantial extrapolation. Calibration (ECE) of the shared parameter  $\mu$  does not clearly differ between the methods.

information processing speed  $\mu$  and associated error  $\sigma$ .

$$\begin{aligned}
 \mu, \sigma, \alpha, \tau, \beta, \eta &\sim p(\mu, \sigma, \alpha, \tau, \beta, \eta) && \text{(prior)} \\
 v_n &\sim \text{Normal}(v | \mu, \sigma) && \text{(per-trial entanglement)} \\
 \mathbf{x}_n &\sim \text{DDM}(\mathbf{x} | \alpha, \tau, v_n, \beta) && \text{(reaction time)} \\
 \mathbf{y}_n &\sim \text{Normal}(\mathbf{y} | v_n, \eta) && \text{(CPP waveform)}
 \end{aligned} \tag{12}$$

In this model,  $\text{DDM}(\mathbf{x} | \cdot)$  denotes the standard (Wiener) drift-diffusion model [50]. Further,  $\text{Normal}(\mathbf{y} | \cdot)$  represents the neurocognitive CPP model from [51]. Crucially, the data sources are entangled on the trial level since the shared information uptake rate  $v_n$  is sampled for each experimental trial  $n$ . This implies an equal number of observations for both sources, corresponding to the number of trials,  $N=M=200$ , per data set  $\mathcal{D}^{(k)} = \{\mathbf{X}^{(k)}, \mathbf{Y}^{(k)}\}$ .

Missing data are a common problem in applied data analyses, and there exists a myriad of methods to tackle missing data in traditional statistical workflows that do not feature

<sup>6</sup>In the cognitive sciences, one trial refers to one event in an overarching experiment (e.g., displaying one image that shall prompt one decision). A whole experiment then consists of tens to thousands of trials [52]

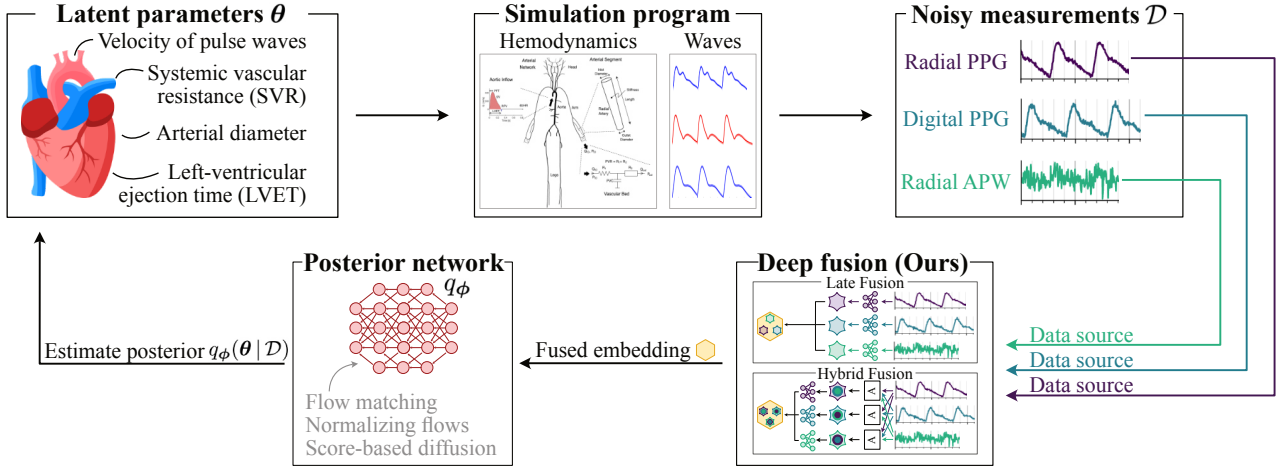


Fig. 6: **Experiment 3.** Overview of the data-generating process (top row) for noisy measurements based on the pulse wave database [53]. These noisy measurements  $\mathcal{D}$  are incompatible due to their varying shape, and our deep fusion techniques integrate the different modalities into a single fused embedding for subsequent neural posterior inference. Illustrations in the box *Simulation program* from [53] under CC-BY 4.0 license.

amortized inference of deep neural networks. We want to study the potential of our deep fusion schemes to handle partially missing data by incorporating missingness into the training phase. To this end, we synthetically induce missing data by uniformly drawing a random missing rate between 1% and 10% for each training batch. Then, we encode missingness with a dedicated ‘missing’ value as well as a binary mask, as proposed by [42].

Notably, in this experiment, it is actually possible to directly concatenate the sources ( $\mathbf{X}, \mathbf{Y}$ ) on the data level because  $\mathbf{X}$  and  $\mathbf{Y}$  have identical matrix shapes. In other words, the neural networks *could* directly process the concatenated data and our fusion schemes are not strictly necessary just due to incompatible data formats. Therefore, we include ‘direct concatenation’ as a baseline for this experiment and observe drastic performance gains from using principled fusion schemes with MultiNPE (see below). With a simulation budget of  $K=4096$ , we compare direct concatenation (baseline), late fusion, and hybrid fusion with respect to the quality of the approximate posterior samples under increasing levels of missing data.

**Results.** As displayed in Figure 5, both late fusion and hybrid fusion outperform direct concatenation via increased accuracy (RMSE) across all missing data rates. The calibration (ECE) on the shared information uptake parameter  $\mu$  does not differ between the methods. This underscores the potential of deep fusion in SBI even in situations where direct concatenation would be possible, and our fusion architectures do not have an advantage by having access to more raw data. We hypothesize that our fusion scheme embodies a favorable inductive bias by separating the data sources through our tailored neural network architectures.

### C. Experiment 3: Cardiovascular model with three sources

Preventing cardiovascular diseases is a fundamental challenge of precision medicine, and scientific hemodynamics simulators are important tools to understand the cardiovascular system [54]. The *pulse wave database* contains data from 4374 *in silico* subjects, and the simulator has previously been validated with *in vivo* data [53]. As illustrated in Figure 6, a whole-body simulator models blood flows through the 116 largest human arteries via a system of differential equations. The output of the simulator are pulse wave measurements of a single heart beat at different locations in the human body, including both photoplethysmograms (PPG) and arterial pressure waveform (APW). In precision medicine, this simulator serves as an *in silico* model that helps researchers study the dynamics of blood flow and associated diseases.

In this experiment, the parameters  $\theta$  and the measurements  $\mathcal{D}$  are only available as a fixed data set. Thus, we do not have access to the prior distribution  $p(\theta)$  or the simulation program  $\theta \mapsto \mathcal{D}$  which defines the forward process from latent parameters  $\theta$  to measurements  $\mathcal{D}$ .<sup>7</sup> This is a particularly challenging case because all parts of the joint probabilistic model  $p(\theta, \mathcal{D})$  are now black-box objects that must be implicitly modeled by our neural network approximators.

Data from the pulse wave database has been previously analyzed with single-source SBI methods [8], and our experiment is closely inspired by this work. As in [8] we aim to use simulation based inference to solve the inverse problem of tracing noisy measurements  $\mathcal{D}$  back to physiological latent parameters  $\theta$  that can explain the measurements. Through our novel addition of attention and fusion mechanisms to learn embeddings for neural posterior estimation, we aim to tackle the additional challenges of using multi-modal cardiovascular

<sup>7</sup>Since we treat the prior and the forward simulation program as black boxes, we do not list the respective mathematical model formulations in this experiment. We refer the interested reader to Charlton et al. [53] for details of the hemodynamics simulator.



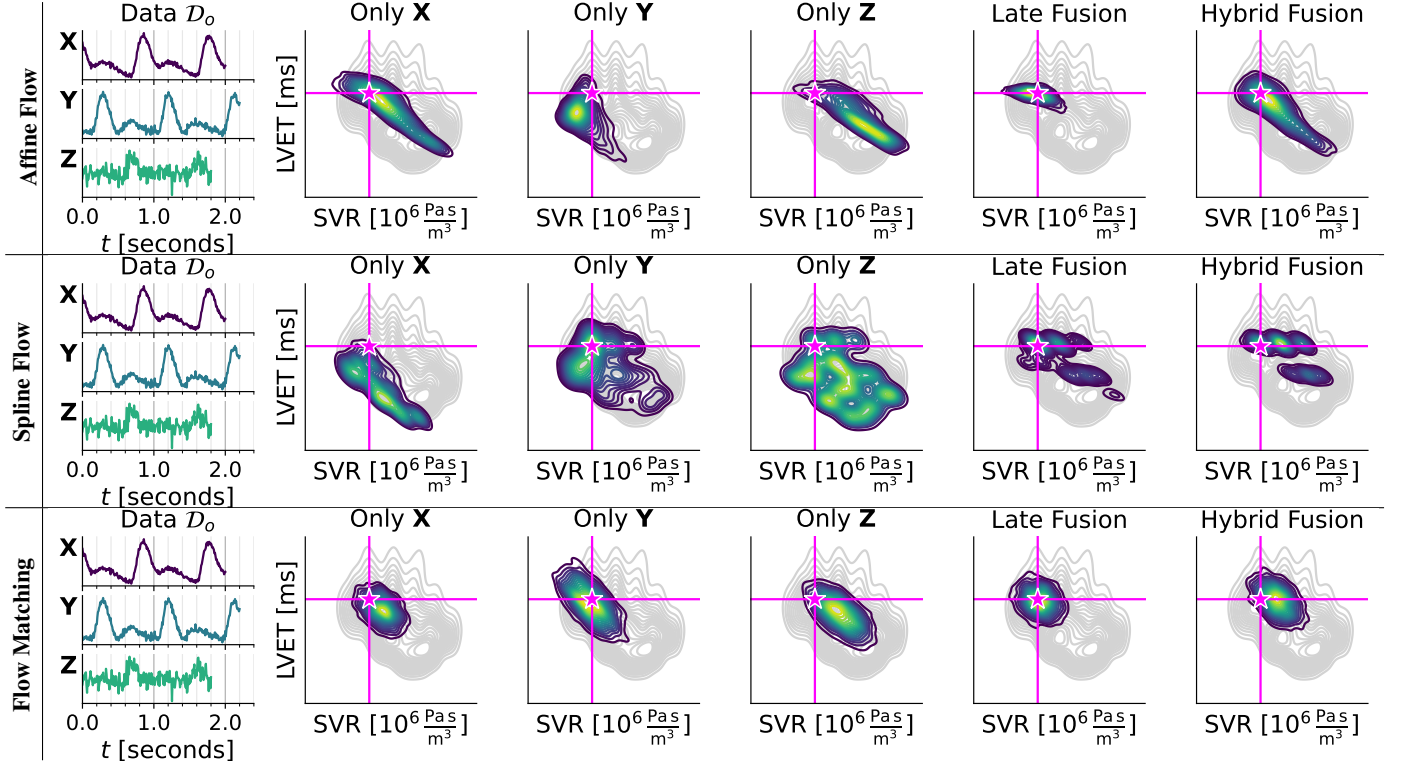


Fig. 7: **Experiment 3:** Bivariate posterior plots of SVR and LVET for one unseen data set  $\mathcal{D}_o$  (left), as obtained by NPE with affine coupling flows, neural spline flows, and flow matching. We show the implicit prior density (gray), the ground-truth  $\theta_*$  (magenta), and a KDE over 10 000 approximate posterior samples (viridis). Visually, affine flows yield the best approximate posterior distributions compared to the other generative model backbones. Within affine flows, posteriors from late fusion and hybrid fusion are more concentrated at the ground-truth compared to the single-source methods, which is generally a desirable property (cf. Table II for calibration results).

data, such as asynchronous measurements and measurement errors that vary between modalities.

We consider the following measurements as individual data sources  $\mathcal{D} = \{\mathbf{X}, \mathbf{Y}, \mathbf{Z}\}$ : PPG at the digital artery ( $\mathbf{X}$ ), PPG at the radial artery ( $\mathbf{Y}$ ), and APW at the radial artery ( $\mathbf{Z}$ ). The shared latent parameters  $\theta$  in this experiment consist of the left ventricular ejection time (LVET), the systemic vascular resistance (SVR), the average diameter of arteries, and the pulse wave velocity.

As an extension to Wehenkel et al. [8], we evaluate the challenging and realistic setting where measurements from different devices ( $\mathbf{X}, \mathbf{Y}, \mathbf{Z}$ ) are not synchronized in time (see Figure 6, top-right panel). Concretely, we achieve this with a two-step process: First, we loop the single-beat signals to a longer sequence and crop the sequence into a fixed-length measurement interval for each subject. The length of the cropped signal differs between the sources ( $\mathbf{X}$ : 2.0 seconds,  $\mathbf{Y}$ : 2.2 seconds,  $\mathbf{Z}$ : 1.8 seconds). The sequence onset times are randomly sampled for each subject and source, which means that the cropped signals are not synchronized anymore within each subject. Second, we add Gaussian white noise to the signals, and the signal-to-noise (SNR) ratio is specific for each data source ( $\mathbf{X}$ : 25dB,  $\mathbf{Y}$ : 20dB,  $\mathbf{Z}$ : 30dB). This emulates different measurement devices in a hospital, where each device has a specific measurement error. In this realistic setting,

we cannot simply concatenate the inputs, but instead require fusion schemes to integrate the heterogeneous cardiovascular measurements  $\mathbf{X}, \mathbf{Y}, \mathbf{Z}$ .

Since the data are only available as an offline data set and not as a simulation program, we face a scenario with both an implicit likelihood and an implicit prior. Further, a subject’s age is a key factor for cardiovascular health in the pulse wave database [53]. Thus, we follow Wehenkel et al. [8] and use the age variable in the data set as an additional direct condition (i.e., without embedding) for the neural density estimator. We compare the performance of single-source models (i.e., only access to  $\mathbf{X}, \mathbf{Y}$ , or  $\mathbf{Z}$ ) to our multimodal variants late fusion and hybrid fusion. Since simulation-based posterior estimation for this data set is a challenging problem by itself, we further employ each strategy with three different generative models: (i) affine flows [55]; (ii) neural spline flows [56]; and (iii) flow matching [57].

**Results.** Overall, affine flows emerge as the best backbone neural network in this experiment, closely followed by flow matching (see Table II). The subpar performance of spline flows may be related to the relatively large data dimension in conjunction with the probabilistic geometry of the noisy time series data. Within both affine flows and flow matching, late fusion shows the best combination of high accuracy and information gain, but suffers from poor calibration. As

Architecture		Time <sup>1</sup> ↓	RMSE ↓	ECE [%] ↓	Contraction ↑
Affine Flow	Only <b>X</b>	1296	0.94	1.90	0.53
	Only <b>Y</b>	1307	1.05	1.38	0.39
	Only <b>Z</b>	1202	0.90	3.47	0.55
	Late Fusion	2010	<b>0.45</b>	5.25	<b>0.90</b>
	Hybrid Fusion	2761	0.75	1.33	0.68
Spline Flow	Only <b>X</b>	1991	1.08	1.27	0.37
	Only <b>Y</b>	2028	1.05	0.82	0.39
	Only <b>Z</b>	1993	1.18	0.97	0.24
	Late Fusion	1880	0.87	<b>0.67</b>	0.56
	Hybrid Fusion	3612	0.64	1.62	0.75
Flow Matching	Only <b>X</b>	1040	0.55	6.93	0.85
	Only <b>Y</b>	1155	0.67	6.29	0.78
	Only <b>Z</b>	<b>993</b>	0.83	8.30	0.72
	Late Fusion	2668	0.50	8.24	0.87
	Hybrid Fusion	4208	0.62	5.98	0.81

TABLE II: **Experiment 3:** Test set performance of different generative models (affine flow, spline flow, flow matching) and fusion strategies (single sources, late fusion, hybrid fusion). Flow matching requires the least time for neural network training on single sources, while it is noticeably slower for the well-performing fusion approaches. We observe that affine flows with late fusion achieve the best (lowest) posterior bias and variance, as quantified by RMSE, and the highest information gain, as evidenced by the highest contraction. However, late fusion leads to a decline in calibration, as indexed by the expected calibration error (ECE), and hybrid fusion alleviates this issue. Taking all three metrics into account, affine flows with late fusion or hybrid fusion yield the best results. <sup>1</sup> Training time [seconds].

opposed to **Experiment 1**, late fusion and hybrid fusion differ with respect to their performance profile: While late fusion yields superior accuracy and contraction, it does not reach the calibration quality of hybrid fusion. Thus we conclude that affine flows with either late fusion or hybrid fusion are desirable for the presented application in precision medicine. [Figure 7](#) shows the bivariate approximate posterior distributions of the parameters SVR and LVET for one observed data set  $\mathcal{D}_o$ . Visual inspection confirms the previously described results, and we recommend affine flows with late or hybrid fusion in this application (see **Supplementary Material** for additional results).

## VI. CONCLUSION

We presented MultiNPE, a new multimodal approach to perform simulation-based Bayesian inference for models with heterogeneous data-generating processes. Our method overcomes the previous inability of neural simulation-based inference algorithms to integrate information from multiple sources. We achieve this information synthesis by constructing expressive embedding architectures which build on state-of-the-art work on deep fusion learning: (i) attention-based early fusion; (ii) late fusion; and (iii) attention-based hybrid fusion. MultiNPE seamlessly combines information from heterogeneous sources, which has previously been infeasible with neural posterior estimators.

We validated MultiNPE on a 10-dimensional benchmark task with a known ground-truth posterior to compare against. Our method showed clearly superior neural network training dynamics, and the resulting posteriors were better than the ones obtained by current state-of-the-art single-source alternatives. In the second experiment, we applied MultiNPE to an applied problem in cognitive neuroscience, where information from brain wave measurements and behavioral response times shall be integrated in a joint integrative model. In this setting, we showed how MultiNPE outperforms the alternative approach even though both algorithms have access to all the data. This effect is particularly pronounced under partially missing observations, which our deep fusion schemes can compensate for. Finally, we showcase how MultiNPE can help medical practitioners integrate information on a patient’s cardiovascular health under realistic settings in a hospital, where measurements are taken with different devices that are not synchronized. This emphasizes the practical utility of our method in applications of precision medicine and real-time health monitoring.

A central research question of this work asked which fusion strategy (i.e., early fusion, late fusion, hybrid fusion) is the most useful for Bayesian inference. Across all experiments, we observed that late fusion and hybrid fusion achieved the best performance, as indexed by parameter recovery (RMSE), probabilistic calibration (ECE), and Bayesian information gain (posterior contraction). For practical applications, we recommend considering both late and hybrid fusion, where the exact choice should be assessed on a case-to-case basis with the probabilistic diagnostics we presented.

Overall, our results underscore the potential of MultiNPE as a novel simulation-based inference tool for real-world problems with multiple data sources. It pushes the boundaries of modern simulation-based inference with neural networks and further paves the way for its widespread application across the sciences and engineering. The **FAQ** section in the Appendix answers some questions we encountered prior to submission.

## Acknowledgments

We thank Lasse Elsemüller for insightful feedback and input on the manuscript. MS thanks the Cyber Valley Research Fund (grant number: CyVy-RF-2021-16), the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) under Germany’s Excellence Strategy EXC-2075 - 390740016 (the Stuttgart Cluster of Excellence SimTech), the Google Cloud Research Credits program, and the European Laboratory for Learning and Intelligent Systems (ELLIS) PhD program for support.

## APPENDIX A

### FREQUENTLY ASKED QUESTIONS (FAQ)

**Q:** Why not just train the summary networks in isolation? Why would you train them jointly with the normalizing flow?

**A:** This is not generally possible because we are interested in (learned) summary statistics that are optimal for *posterior inference*. The summaries are not meant to reconstruct the

data. Thus, it is paramount that the summary network(s) and the neural density estimator are trained end-to-end.

**Q:** Why did you choose these tasks? Is there no benchmark suite for multimodal simulation-based inference?

**A:** Since this paper is pioneering the joint analysis of heterogeneous data sources in simulation-based inference, there is no established benchmark suite to use. We aimed to cover a wide range of practically relevant cases with our selected experiments, ranging from a toy example (**Experiment 1**), to a practically important missing data setting on a recent cognitive model (**Experiment 2**), to in-silico cardiovascular data with realistic measurement quality challenges (**Experiment 3**).

**Q:** Why not benchmark against multimodal SBI methods?

**A:** See above, our work is the first SBI method for fusing heterogeneous data. We compare our approach with single-source methods in **Experiments 1 and 3**, and we specifically designed **Experiment 2** such that the naïve SBI approach can handle the two-source input if the data are concatenated.

**Q:** Can the method be applied to other multimodal data sets like text and images?

**A:** While it is theoretically possible for our deep fusion schemes, the simulation-based inference approach with full uncertainty quantification will likely not scale to such data at the frontier of generative AI research.

## APPENDIX B IMPLEMENTATION DETAILS

All experiments are performed on a machine with 4 vCPUs, an NVIDIA T4 GPU, and 15GB RAM.

### A. Experiment 1

**Neural network details** All transformer embedding networks use 4 attention heads, 32-dimensional keys, 10% dropout, layer normalization, 2 fully-connected layers of 64 units each within the attention heads, and learn 10-dimensional embeddings. The multihead-attention blocks for early fusion and the early stage of hybrid fusion use 4 attention heads, 32-dimensional keys, 10% dropout, layer normalization, and 3 fully-connected layers of 64 units each within the attention heads. The conditional normalizing flow consists of 8 affine coupling layers, each with one fully-connected layer of 32 units and an L2 kernel regularizer with weight  $\gamma = 10^{-4}$ . Across all architectures, we use an initial learning rate of  $10^{-4}$  with cosine decay, a batch size of 32, and train for 30 epochs without early stopping.

### B. Experiment 2

The prior distributions for the parameters are defined as

$$\begin{aligned}\mu &\sim \mathcal{U}(0.1, 3), \\ \alpha &\sim \mathcal{U}(0.5, 2), \\ \beta &\sim \mathcal{U}(0.1, 0.9), \\ \tau &\sim \mathcal{U}(0.1, 1), \\ \sigma &\sim \mathcal{U}(0, 2), \\ \eta &\sim \mathcal{U}(0, 2),\end{aligned}\tag{13}$$

where  $\mathcal{U}(a, b)$  denotes the uniform distribution with lower bound  $a$  and upper bound  $b$ .

**Missing data** We synthetically induce missing data in the data generating process by uniformly sampling individual missing rates  $\rho_x, \rho_y \in [0.01, 0.10]^2$  for each batch during training. Subsequently, we create two independent masks  $m_x \sim \text{Bernoulli}(1 - \rho_x)$  and  $m_y \sim \text{Bernoulli}(1 - \rho_y)$  which determine whether each data set is missing or not. As proposed by [42] for simulation-based inference, we encode missing data as a constant  $\mathbf{c}$  with measure zero under the data generating process,  $\mathbf{c} = -1.0, p(\mathbf{c}) = 0$ . Additionally, we append the masks  $m_x, m_y$  to the data  $\mathbf{X}, \mathbf{Y}$ , which has been shown to facilitate discrimination between available and missing data for the neural density estimator [42]. **Neural network details**

The multi-head attention blocks for early fusion and hybrid fusion use no layer normalization. The embedding networks are equivariant set transformers with 2 self-attention blocks of 4 attention heads, 64-dimensional keys and 10% dropout. The number of learned embeddings equals 12, which implements the heuristic from [58] to use twice the number of inference target parameters. The neural density estimator is a neural spline flow with 4 coupling blocks, each consisting of 3 dense blocks with 128 units, L2 kernel regularization with weight  $\gamma = 10^{-4}$ , 10% dropout, and spectral normalization to further support learning in low data regimes with missing data. All networks train for 100 epochs with a batch size of 32.

### C. Experiment 3

In addition to the preprocessing steps outlined in the main text, we downsample the original 500Hz signal from the pulsewave database to 125Hz with the naïve method of using only every 4<sup>th</sup> measurement. Further, we normalize data and parameters with respect to the empirical mean and standard deviation of the training set. In order to employ temporal fusion transformers as summary networks, we add a linear time encoding to the data, which is not synchronized between sources.

**Neural network details** The multi-head attention blocks in early fusion and hybrid fusion use 4 attention heads, 32-dimensional keys, 1% dropout, and layer normalization. The temporal fusion transformers that we employ as summary networks for time series data use 2 multi-head self-attention blocks with 4 attention heads each, 32-dimensional keys, 10% dropout, and layer normalization. The embedding networks learn 30-dimensional representations (aka. summary statistics or features). Both the affine coupling flow and the neural spline

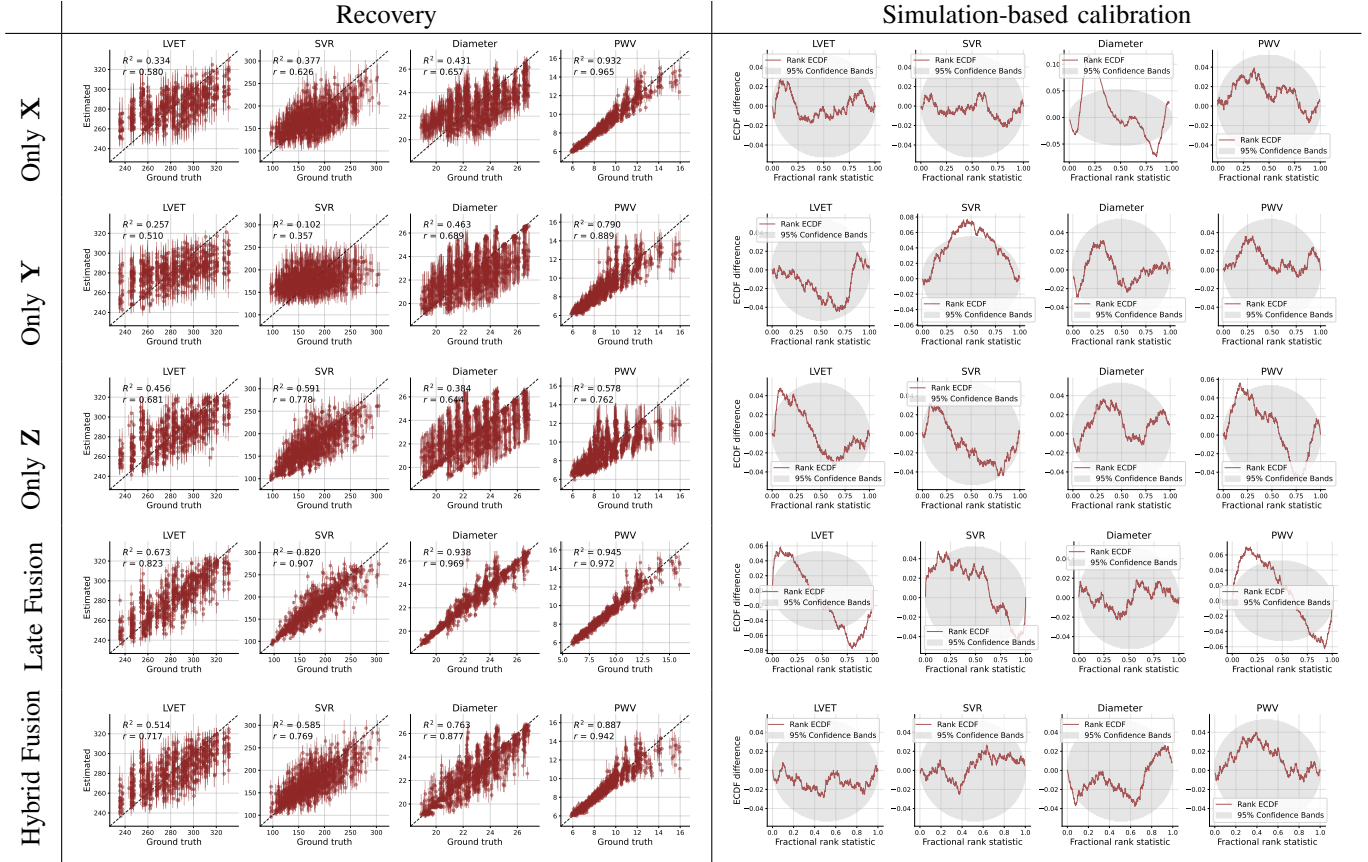
flow use 8 coupling blocks, each featuring 2 dense layers of 64 units, 10% dropout, and an L2 kernel regularizer with weight  $\gamma = 10^{-4}$ . Flow matching uses a drift network with two dense layers of 64 units, 10% dropout and an L2 kernel regularizer with weight  $\gamma = 10^{-4}$ . All architectures use a batch size of 16 and an initial learning rate of  $5 \cdot 10^{-4}$  with subsequent cosine decay. We train the affine coupling flow and neural spline flow for 100 epochs, while flow matching trains for 200 epochs.

**Additional detailed results** In addition to the bivariate posterior plots in the main text, we show results for further test instances and all three neural density estimators. Additionally, we report additional results on the closed-world performance over the entire test set, namely the (i) parameter recovery (ground-truth vs. estimated); and (ii) detailed simulation-based calibration (SBC) analyses (see Figure 8, 9, 10).

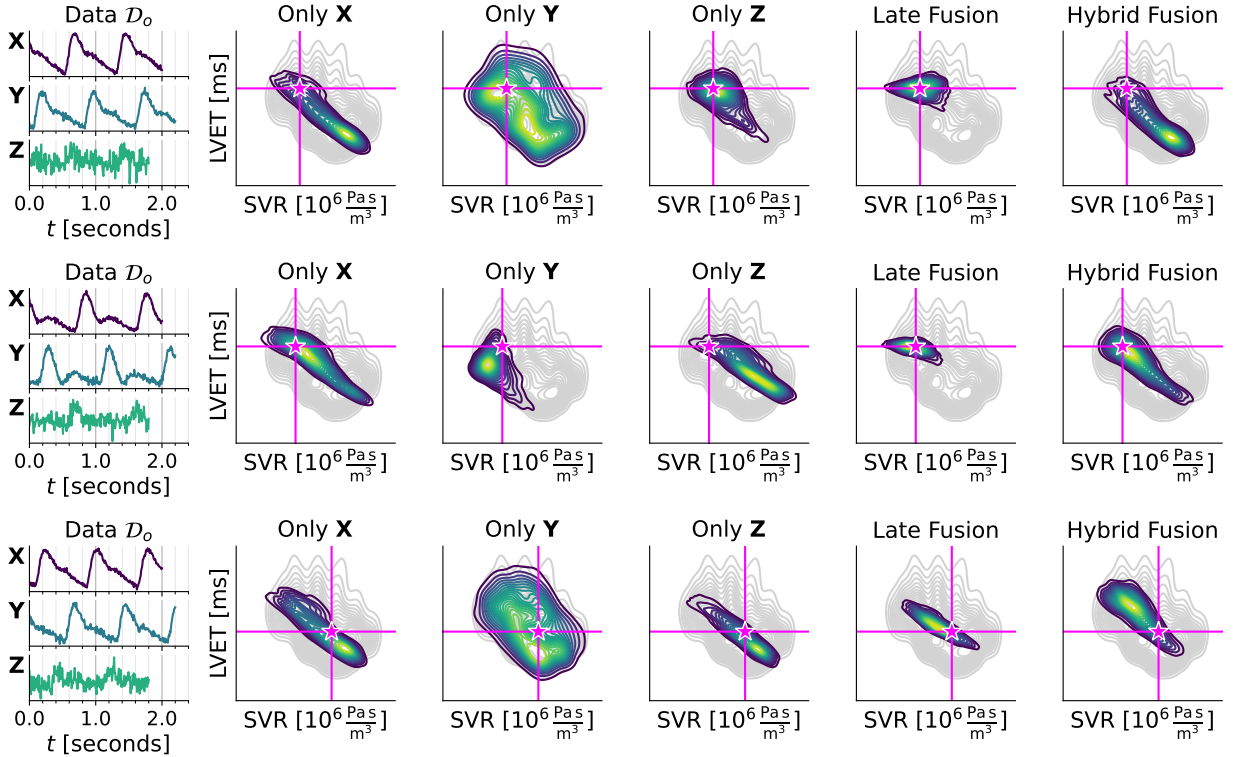
## REFERENCES

- [1] A. Lavin, H. Zenil, B. Paige *et al.*, “Simulation intelligence: Towards a new generation of scientific methods,” *arXiv preprint*, 2021.
- [2] P. J. Green, K. Łatuszyński, M. Pereyra, and C. P. Robert, “Bayesian computation: a summary of the current state, and samples backwards and forwards,” *Statistics and Computing*, vol. 25, no. 4, p. 835–862, Jun. 2015.
- [3] E. Štrumbelj, A. Bouchard-Côté, J. Corander, A. Gelman, H. Rue, L. Murray, H. Pesonen, M. Plummer, and A. Vehtari, “Past, present and future of software for bayesian inference,” *Statistical Science*, vol. 39, no. 1, Feb. 2024.
- [4] S. A. Sisson, Y. Fan, and M. Beaumont, *Handbook of approximate Bayesian computation*. CRC Press, 2018.
- [5] J.-M. Marin, P. Pudlo, C. P. Robert, and R. J. Ryder, “Approximate Bayesian computational methods,” *Statistics and Computing*, 2012.
- [6] R. M. Neal, *MCMC using Hamiltonian dynamics*, May 2011.
- [7] D. M. Blei, A. Kucukelbir, and J. D. McAuliffe, “Variational inference: A review for statisticians,” *Journal of the American statistical Association*, vol. 112, no. 518, pp. 859–877, 2017.
- [8] A. Wehenkel, J. Behrmann, A. C. Miller, G. Sapiro, O. Sener, M. Cuturi, and J.-H. Jacobsen, “Simulation-based inference for cardiovascular models,” 2023.
- [9] G. Papamakarios, E. Nalisnick, D. J. Rezende, S. Mohamed, and B. Lakshminarayanan, “Normalizing flows for probabilistic modeling and inference,” *J. Mach. Learn. Res.*, vol. 22, no. 1, jan 2021.
- [10] L. Sharrock, J. Simons, S. Liu, and M. Beaumont, “Sequential neural score estimation: Likelihood-free inference with conditional score based diffusion models,” in *Proceedings of the 41st International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, R. Salakhutdinov, Z. Kolter, K. Heller, A. Weller, N. Oliver, J. Scarlett, and F. Berkenkamp, Eds., vol. 235. PMLR, 21–27 Jul 2024, pp. 44565–44602.
- [11] M. Dax, J. Wildberger, S. Buchholz, S. R. Green, J. H. Macke, and B. Schölkopf, “Flow matching for scalable simulation-based inference,” in *Neural Information Processing Systems*, 2023.
- [12] M. Schmitt, V. Pratz, U. Köthe, P.-C. Bürkner, and S. T. Radev, “Consistency models for scalable and fast simulation-based inference,” 2023, arXiv:2312.05440.
- [13] L. Raynal, J.-M. Marin, P. Pudlo, M. Ribatet, C. P. Robert, and A. Estoup, “Abc random forests for bayesian parameter inference,” *Bioinformatics*, vol. 35, no. 10, pp. 1720–1728, 2019.
- [14] J. J. Palestro, G. Bahg, P. B. Sederberg, Z.-L. Lu, M. Steyvers, and B. M. Turner, “A tutorial on joint models of neural and behavioral measures of cognition,” *Journal of Mathematical Psychology*, vol. 84, pp. 20–48, 2018.
- [15] M. G. Blum, M. A. Nunes, D. Prangle, and S. A. Sisson, “A comparative review of dimension reduction methods in approximate bayesian computation,” *Statistical Science*, vol. 28, no. 2, pp. 189–208, 2013.
- [16] P. Fearnhead and D. Prangle, “Constructing summary statistics for approximate bayesian computation: semi-automatic approximate bayesian computation,” *Journal of the Royal Statistical Society Series B: Statistical Methodology*, vol. 74, no. 3, pp. 419–474, 2012.
- [17] S. T. Radev, U. K. Mertens, A. Voss, L. Ardizzone, and U. Köthe, “BayesFlow: Learning complex stochastic models with invertible neural networks,” *IEEE transactions on neural networks and learning systems*, 2020.
- [18] Y. Chen, D. Zhang, M. U. Gutmann, A. Courville, and Z. Zhu, “Neural approximate sufficient statistics for implicit models,” in *International Conference on Learning Representations*, 2021.
- [19] J. Chan, V. Perrone, J. Spence, P. Jenkins, S. Mathieson, and Y. Song, “A likelihood-free inference framework for population genetic data using exchangeable neural networks,” *Advances in neural information processing systems*, vol. 31, 2018.
- [20] D. Huang, A. Bharti, A. H. Souza, L. Acerbi, and S. Kaski, “Learning robust statistics for simulation-based inference under model misspecification,” in *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- [21] M. Zaheer, S. Kottur, S. Ravanbakhsh, B. Poczos, R. Salakhutdinov, and A. Smola, “Deep sets,” 2017.
- [22] J. Lee, Y. Lee, J. Kim, A. Kosiorek, S. Choi, and Y. W. Teh, “Set transformer: A framework for attention-based permutation-invariant neural networks,” in *Proceedings of the 36th International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, K. Chaudhuri and R. Salakhutdinov, Eds., vol. 97. US: PMLR, 09–15 Jun 2019, pp. 3744–3753.
- [23] S. Hochreiter and J. Schmidhuber, “Long short-term memory,” *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [24] Q. Wen, T. Zhou, C. Zhang, W. Chen, Z. Ma, J. Yan, and L. Sun, “Transformers in time series: A survey,” in *International Joint Conference on Artificial Intelligence (IJCAI)*, 2023.
- [25] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit,



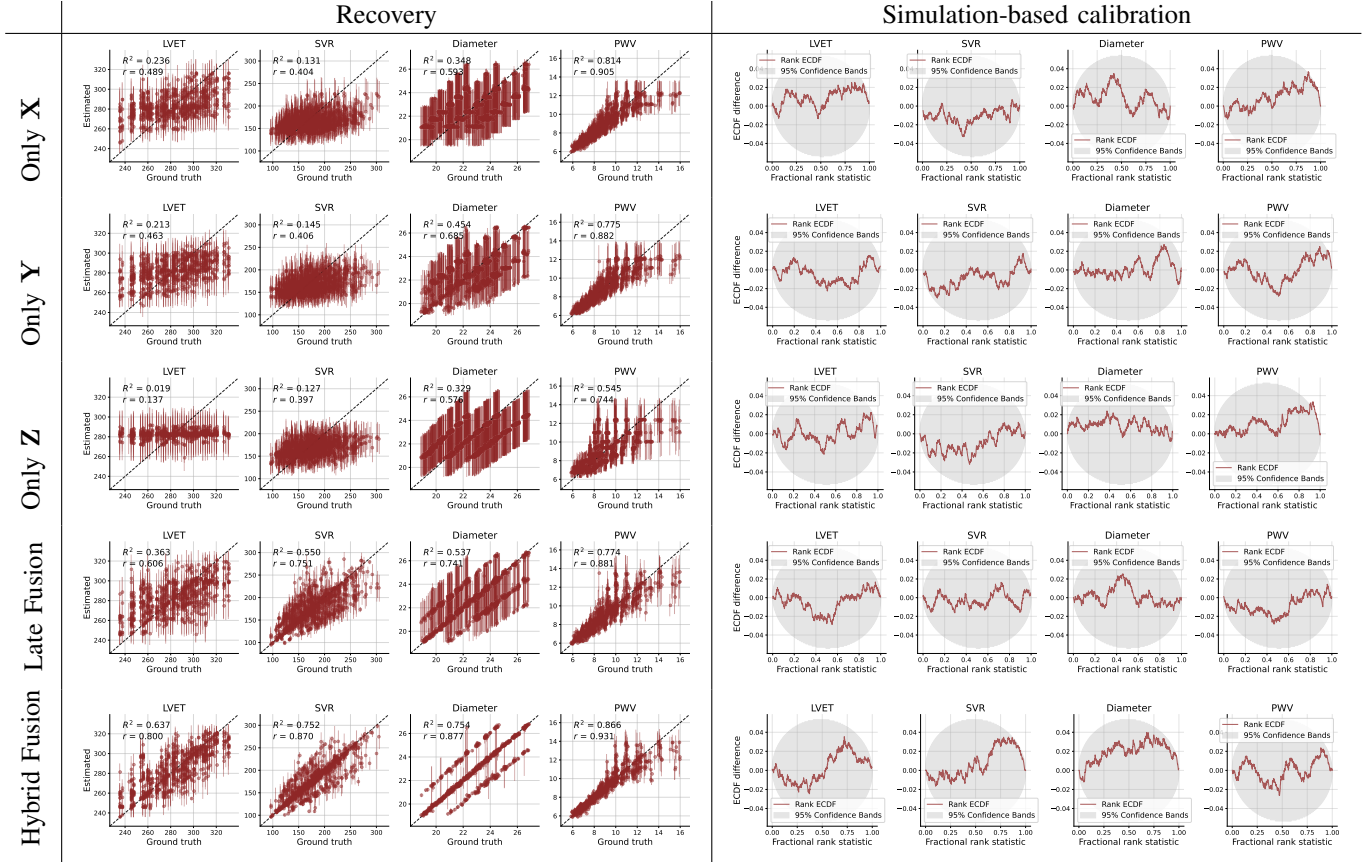


(a) Recovery and simulation-based calibration across the test set

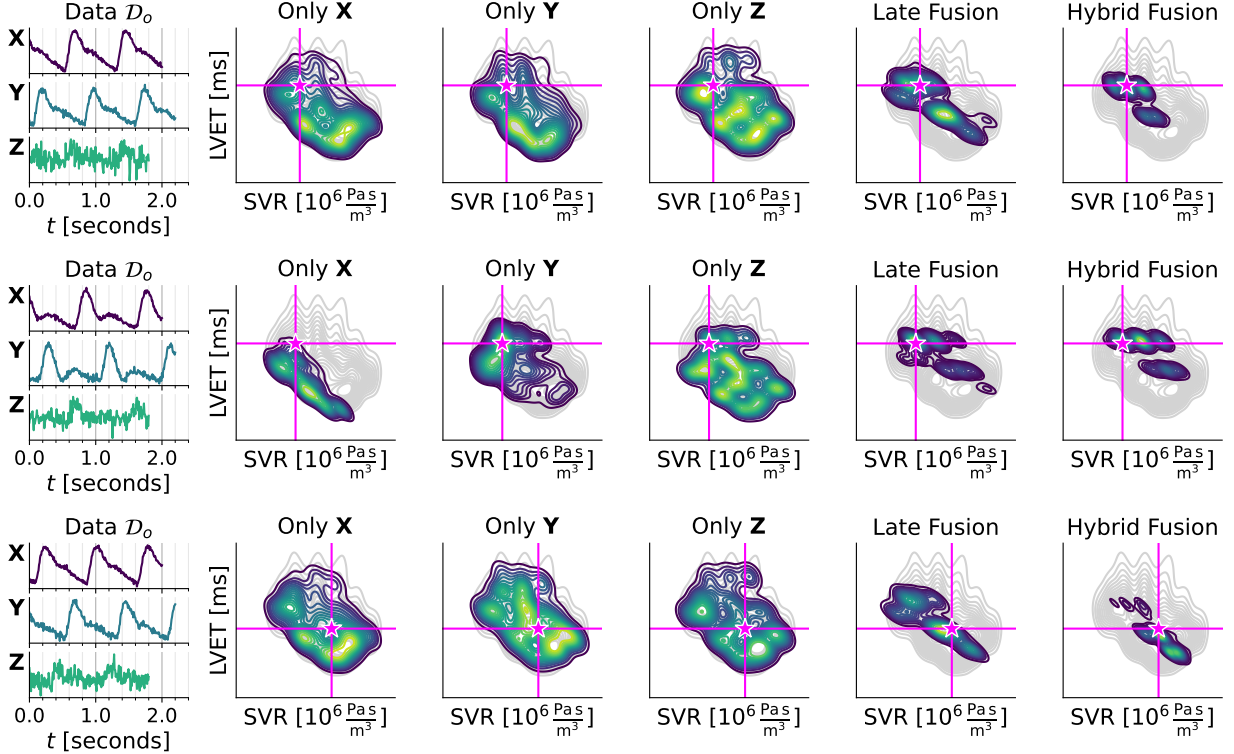


(b) Bivariate posterior plots on further test instances

Fig. 8: Experiment 3, affine coupling flow.

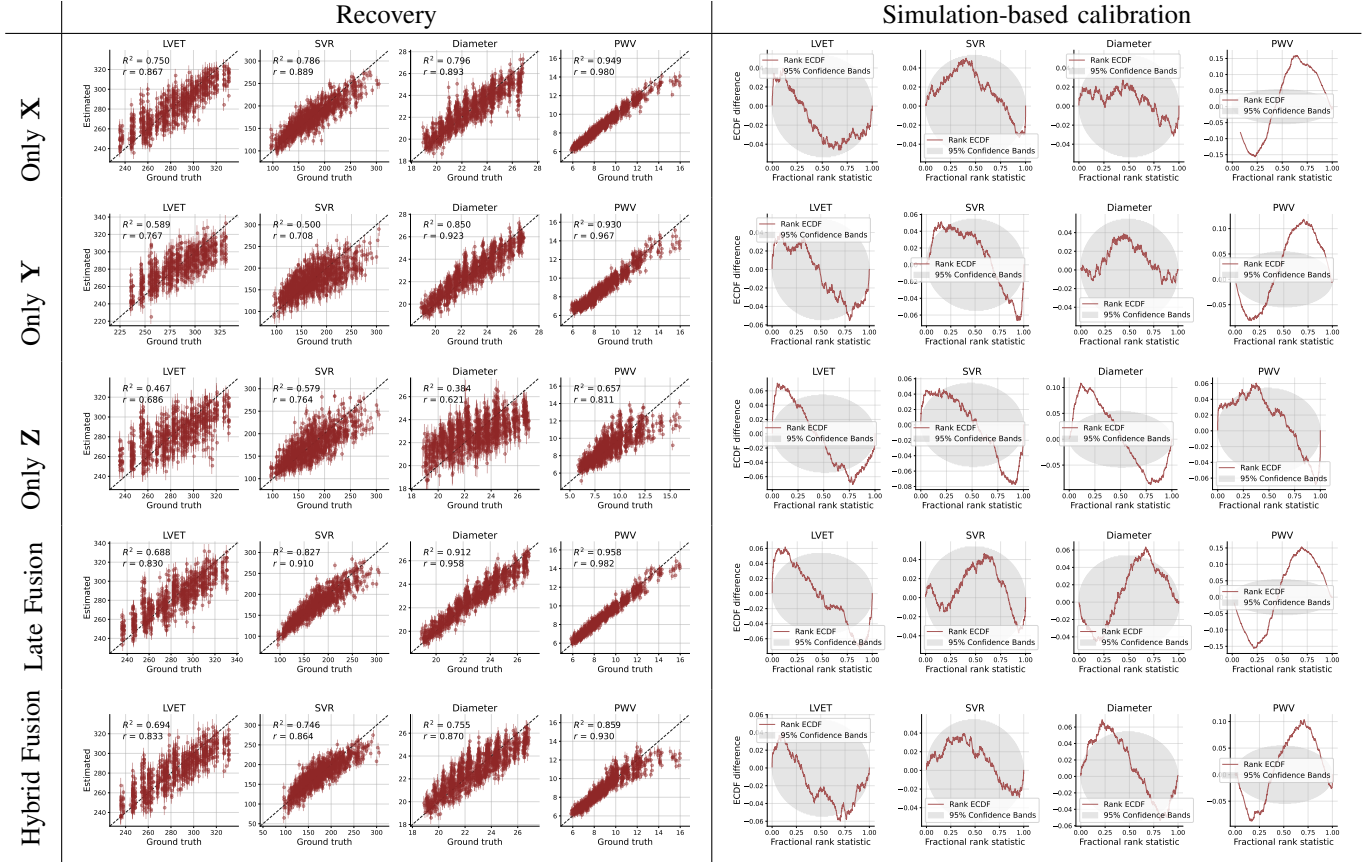


(a) Recovery and simulation-based calibration across the test set

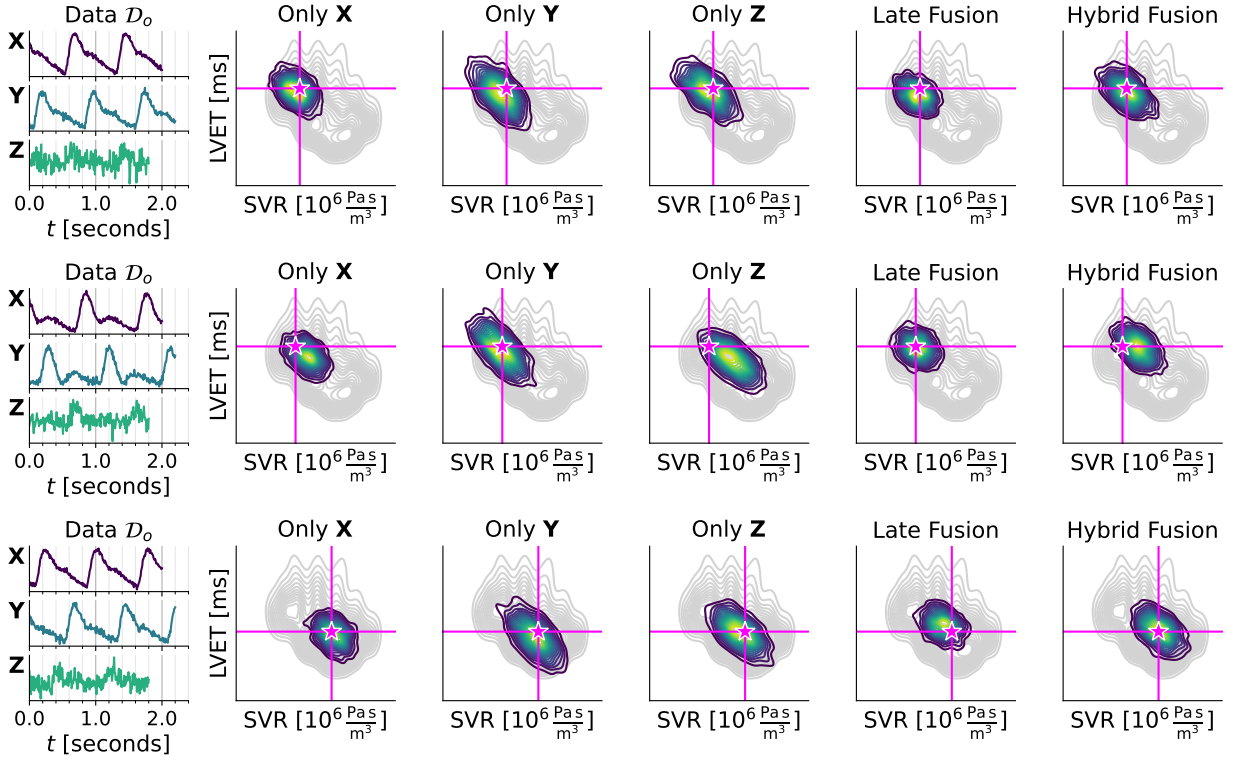


(b) Bivariate posterior plots on further test instances

Fig. 9: Experiment 3, neural spline flow.



(a) Recovery and simulation-based calibration across the test set



(b) Bivariate posterior plots on further test instances

Fig. 10: Experiment 3, flow matching.

- L. Jones, A. N. Gomez, Łukasz Kaiser, and I. Polosukhin, “Attention is all you need,” in *Advances in Neural Information Processing Systems*, I. Guyon, U. von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, Eds., vol. 30, 2017.
- [26] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby, “An image is worth 16x16 words: Transformers for image recognition at scale,” in *International Conference on Learning Representations*, 2021.
- [27] P. K. Atrey, M. A. Hossain, A. E. Saddik, and M. S. Kankanhalli, “Multimodal fusion for multimedia analysis: a survey,” *Multimedia Systems*, vol. 16, no. 6, pp. 345–379, Apr. 2010.
- [28] H. Gunes and M. Piccardi, “Affect recognition from face and body: Early fusion vs. late fusion,” in *2005 IEEE International Conference on Systems, Man and Cybernetics*. US: IEEE, 2005.
- [29] C. Zhang, Z. Yang, X. He, and L. Deng, “Multimodal intelligence: Representation learning, information fusion, and applications,” 2019.
- [30] J. Lu, D. Batra, D. Parikh, and S. Lee, *ViLBERT: Pre-training Task-Agnostic Visiolinguistic Representations for Vision-and-Language Tasks*. Red Hook, NY, USA: Curran Associates Inc., 2019.
- [31] V. Murahari, D. Batra, D. Parikh, and A. Das, “Large-scale pretraining for visual dialog: A simple state-of-the-art baseline,” in *Computer Vision – ECCV 2020*, A. Vedaldi, H. Bischof, T. Brox, and J.-M. Frahm, Eds. Cham: Springer International Publishing, 2020, pp. 336–352.
- [32] J. Ngiam, A. Khosla, M. Kim, J. Nam, H. Lee, and A. Y. Ng, “Multimodal deep learning,” in *Proceedings of the 28th International Conference on International Conference on Machine Learning*, ser. ICML’11. Madison, WI, USA: Omnipress, 2011, p. 689–696.
- [33] W. Guo, J. Wang, and S. Wang, “Deep multimodal representation learning: A survey,” *IEEE Access*, vol. 7, pp. 63 373–63 394, 2019.
- [34] P. Xu, X. Zhu, and D. A. Clifton, “Multimodal learning with transformers: A survey,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 45, no. 10, pp. 12 113–12 132, oct 2023.
- [35] X. Wei, T. Zhang, Y. Li, Y. Zhang, and F. Wu, “Multimodality cross attention network for image and sentence matching,” in *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. US: IEEE, Jun. 2020.
- [36] F. Qingyun, H. Dapeng, and W. Zhaokui, “Cross-modality fusion transformer for multispectral object detection,” 2021.
- [37] K. Gavriluk, R. Sanford, M. Javan, and C. G. M. Snoek, “Actor-transformers for group activity recognition,” in *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 836–845.
- [38] T. Zhi-Xuan, H. Soh, and D. C. Ong, “Factorized inference in deep markov models for incomplete multimodal time series,” 2019.
- [39] P. P. Liang, Z. Liu, Y.-H. H. Tsai, Q. Zhao, R. Salakhutdinov, and L.-P. Morency, “Learning representations from imperfect time series data via tensor rank regularization,” 2019.
- [40] L. Fei-Fei and P. Perona, “A bayesian hierarchical model for learning natural scene categories,” in *2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR’05)*, vol. 2, 2005, pp. 524–531 vol. 2.
- [41] M. Ma, J. Ren, L. Zhao, S. Tulyakov, C. Wu, and X. Peng, “Smil: Multimodal learning with severely missing modality,” 2021.
- [42] Z. Wang, J. Hasenauer, and Y. Schälte, “Missing data in amortized simulation-based neural posterior estimation,” 2023.
- [43] A. Gelman, J. B. Carlin, H. S. Stern, D. B. Dunson, A. Vehtari, and D. B. Rubin, *Bayesian Data Analysis (3rd Edition)*. CRC: Chapman and Hall, 2013.
- [44] D. Habermann, M. Schmitt, L. Kühmichel, A. Bulling, S. T. Radev, and P.-C. Bürkner, “Amortized bayesian multilevel models,” 2024, arXiv:2408.13230.
- [45] C. K. Winkle, “Hierarchical bayesian models for predicting the spread of ecological processes,” *Ecology*, vol. 84, no. 6, pp. 1382–1394, 2003.
- [46] B. Carpenter, A. Gelman, M. D. Hoffman, D. Lee, B. Goodrich, M. Betancourt, M. Brubaker, J. Guo, P. Li, and A. Riddell, “Stan: A probabilistic programming language,” *Journal of statistical software*, vol. 76, no. 1, 2017.
- [47] S. Talts, M. Betancourt, D. Simpson, A. Vehtari, and A. Gelman, “Validating Bayesian inference algorithms with simulation-based calibration,” *arXiv preprint*, 2018.
- [48] P.-C. Bürkner, M. Scholz, and S. T. Radev, “Some models are useful, but how do we know which ones? towards a unified bayesian model taxonomy,” 2022.
- [49] A. Gretton, K. Borgwardt, M. Rasch, B. Schölkopf, and A. Smola, “A Kernel Two-Sample Test,” *The Journal of Machine Learning Research*, vol. 13, pp. 723–773, 2012.
- [50] A. Voss, K. Rothermund, and J. Voss, “Interpreting the parameters of the diffusion model: An empirical validation,” *Memory & Cognition*, vol. 32, no. 7, pp. 1206–1220, 2004.
- [51] A. Ghaderi-Kangavari, J. A. Rad, and M. D. Nunez, “A general integrative neurocognitive modeling framework to jointly describe EEG and decision-making on single trials,” *Computational Brain and Behavior*, 2023.
- [52] V. Lerche, A. Voss, and M. Nagler, “How many trials are required for parameter estimation in diffusion modeling? a comparison of different optimization criteria,” *Behavior Research Methods*, vol. 49, no. 2, p. 513–537, Jun. 2016.
- [53] P. H. Charlton, J. M. Harana, S. Vennin, Y. Li, P. Chowienczyk, and J. Alastruey, “Modeling arterial pulse waves in healthy aging: a database for in silico evaluation of hemodynamics and pulse wave indexes,” *American Journal of Physiology-Heart and Circulatory Physiology*, vol. 317, no. 5, pp. H1062–H1085, 2019.
- [54] E. A. Ashley, “Towards precision medicine,” *Nature*



- Reviews Genetics*, vol. 17, no. 9, pp. 507–522, 2016.
- [55] L. Dinh, J. Sohl-Dickstein, and S. Bengio, “Density estimation using real NVP,” *arXiv preprint arXiv:1605.08803*, 2016.
- [56] C. Durkan, A. Bekasov, I. Murray, and G. Papamakarios, “Neural spline flows,” *Advances in neural information processing systems*, vol. 32, 2019.
- [57] X. Liu, C. Gong, and Q. Liu, “Flow straight and fast: Learning to generate and transfer data with rectified flow,” 2022.
- [58] M. Schmitt, P.-C. Bürkner, U. Köthe, and S. T. Radev, “Detecting Model Misspecification in Amortized Bayesian Inference with Neural Networks,” in *Pattern Recognition. DAGM GCPR 2023. Lecture Notes in Computer Science*. Switzerland: Springer Nature, 2024, p. 541–557.

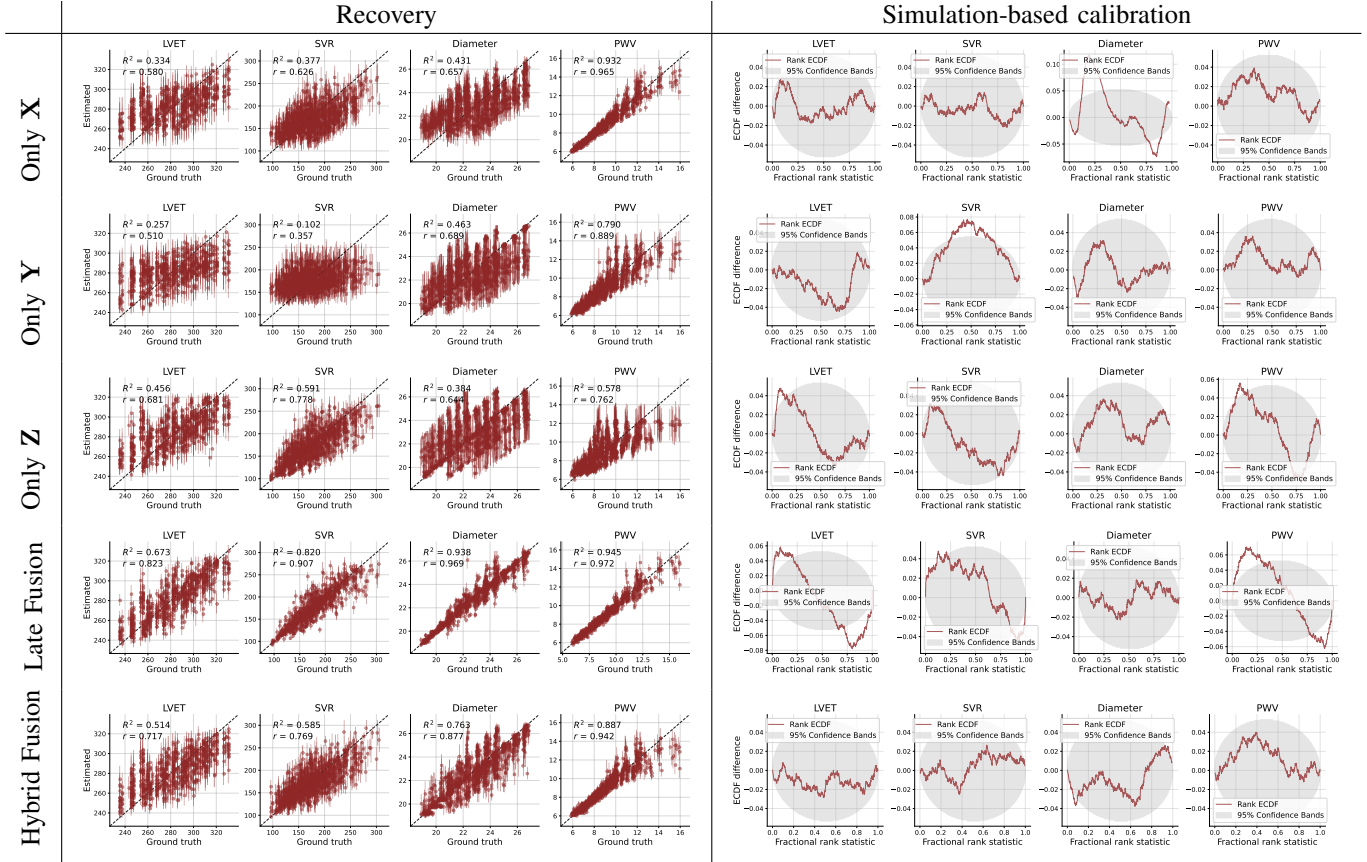
# Fuse It or Lose It: Deep Fusion for Multimodal Simulation-Based Inference (Supplementary Material)

Anonymous Authors

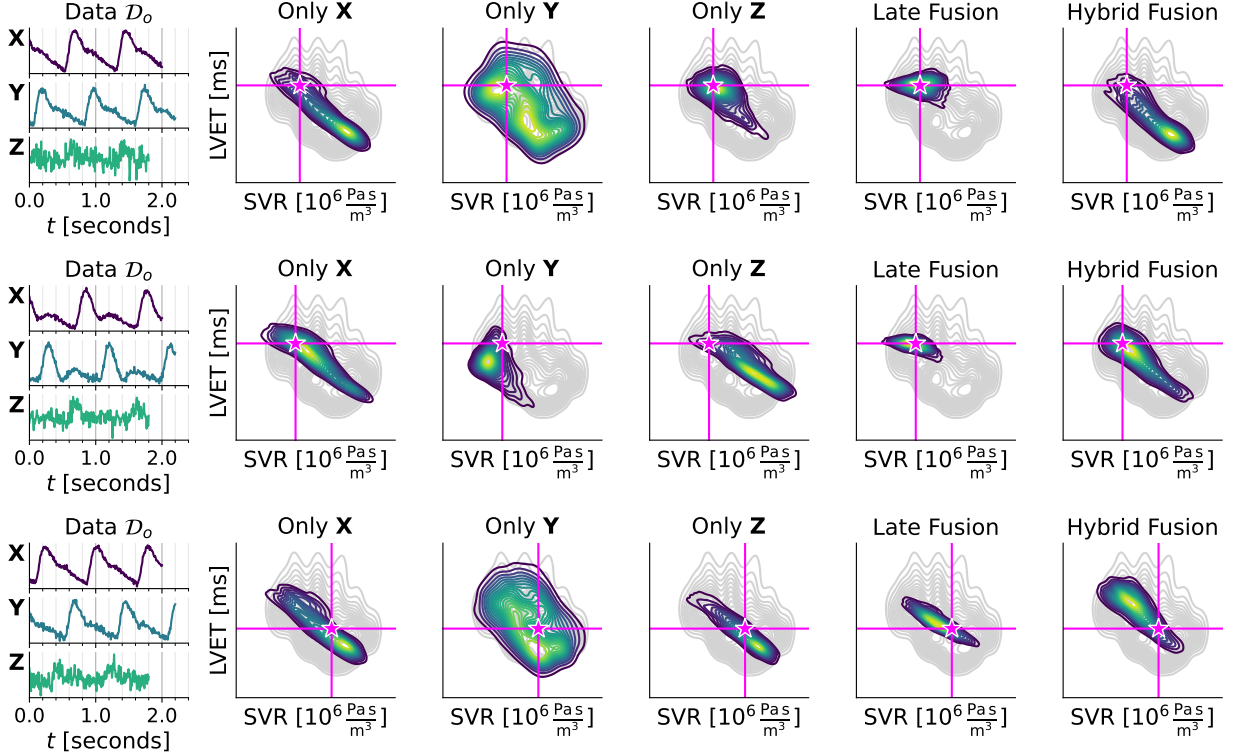
Anonymous Affiliation

## EXPERIMENT 3: ADDITIONAL DETAILED RESULTS

In addition to the bivariate posterior plots in the main text, we show results for further test instances and all three neural density estimators. Additionally, we report additional results on the closed-world performance over the entire test set, namely the (i) parameter recovery (ground-truth vs. estimated); and (ii) detailed simulation-based calibration (SBC) analyses (see [Figure 1, 2, 3](#)).

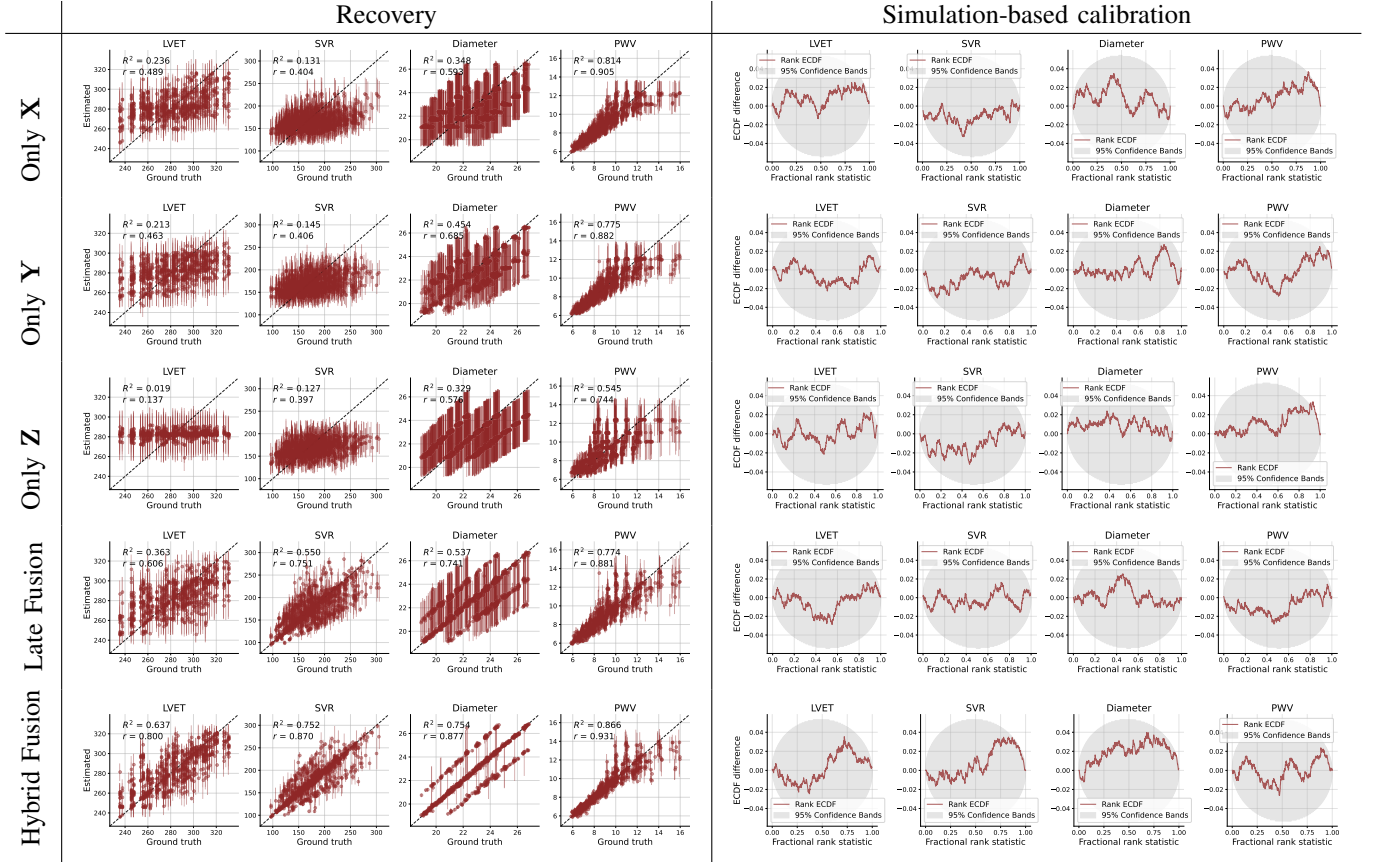


(a) Recovery and simulation-based calibration across the test set

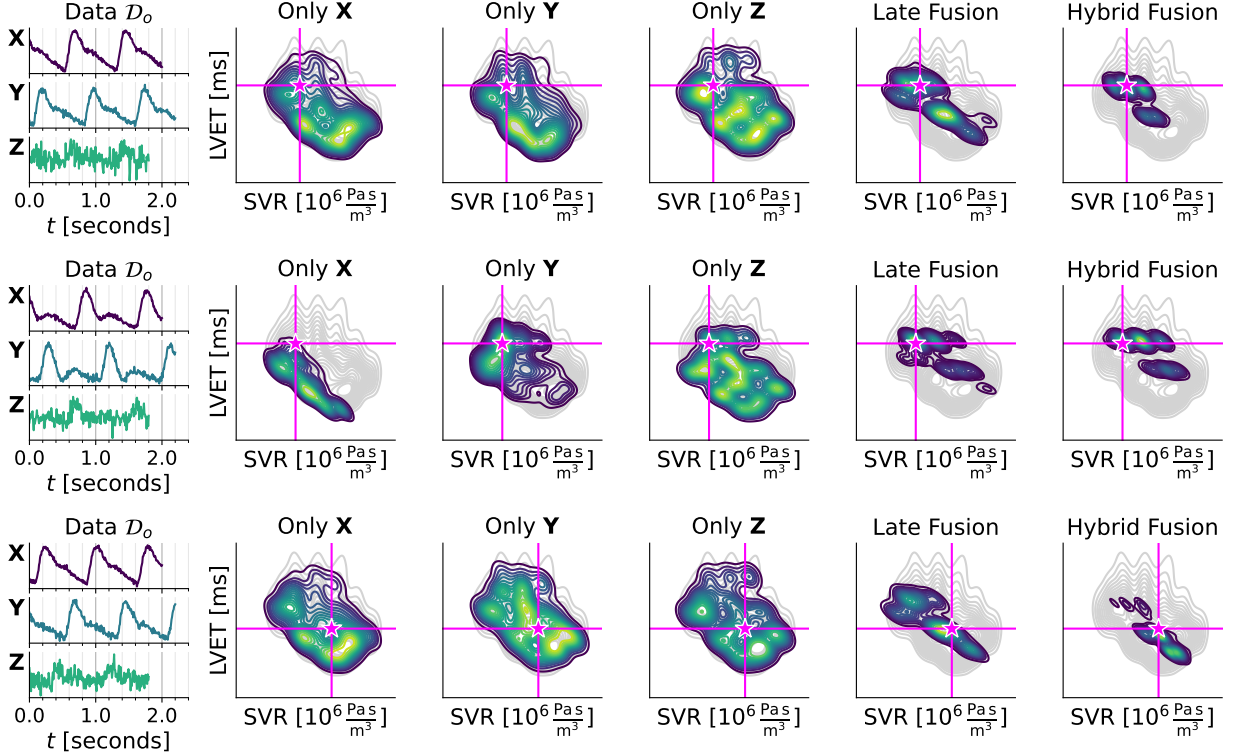


(b) Bivariate posterior plots on further test instances

Fig. 1: Experiment 3, affine coupling flow.



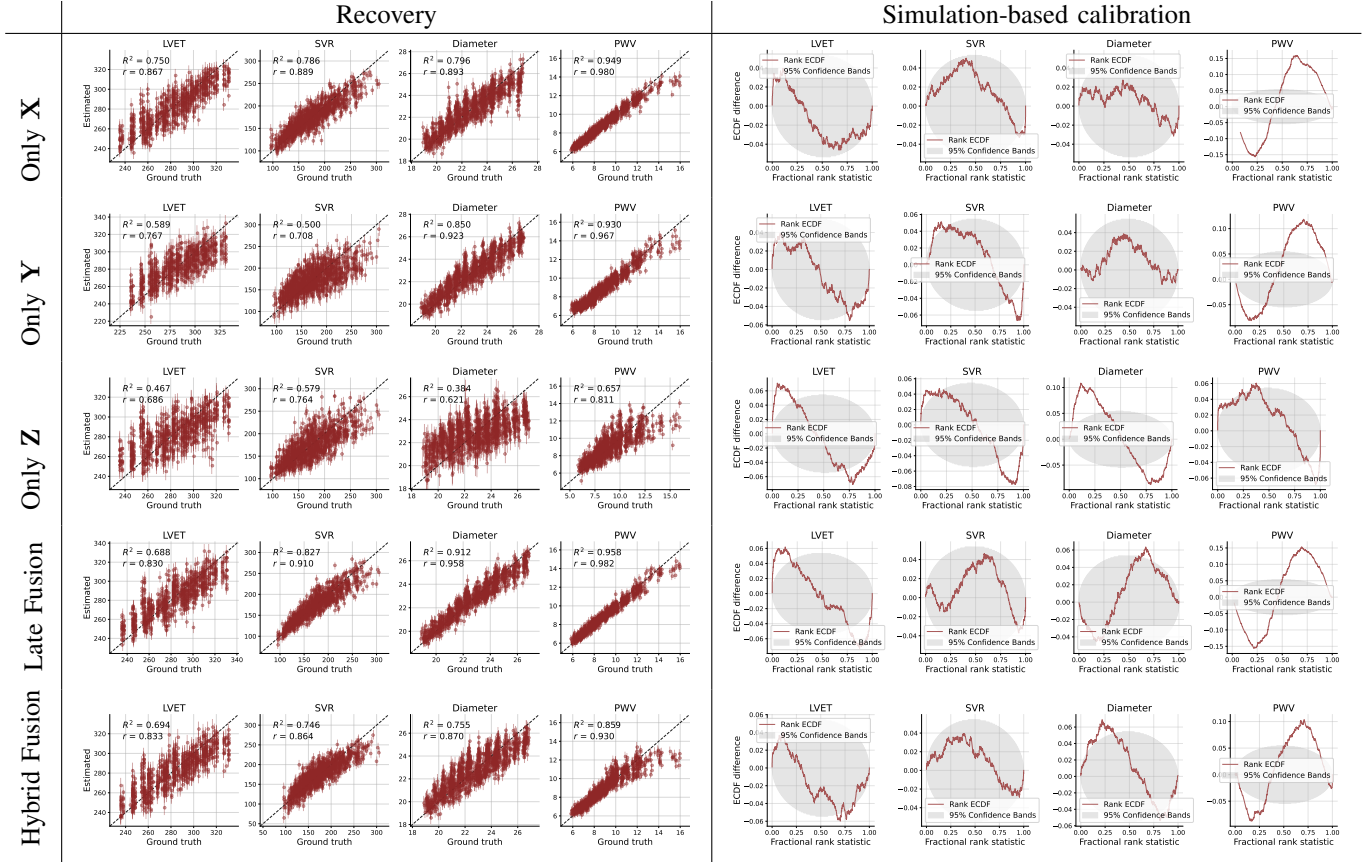
(a) Recovery and simulation-based calibration across the test set



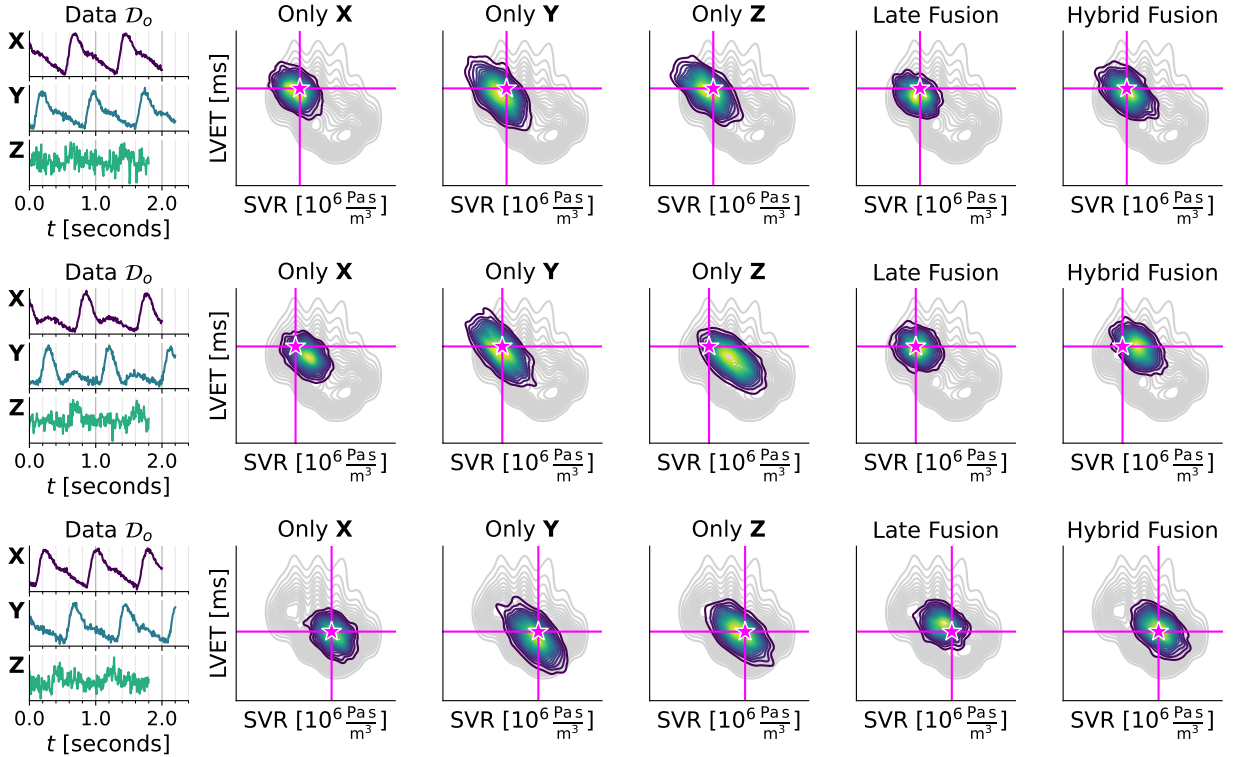
(b) Bivariate posterior plots on further test instances

Fig. 2: Experiment 3, neural spline flow.





(a) Recovery and simulation-based calibration across the test set



(b) Bivariate posterior plots on further test instances

Fig. 3: Experiment 3, flow matching.