

# Weakly Supervised Bayesian Shape Modeling from Unsegmented Medical Images

Jadie Adams<sup>1,2</sup>, Krithika Iyer<sup>1,2</sup>, and Shireen Y. Elhabian<sup>1,2</sup>

<sup>1</sup> Scientific Computing and Imaging Institute, University of Utah, UT, USA

<sup>2</sup> Kahlert School of Computing, University of Utah, UT, USA  
{jadie, iyerkrithika, shireen}@sci.utah.edu

**Abstract.** Anatomical shape analysis plays a pivotal role in clinical research and hypothesis testing, where the relationship between form and function is paramount. Correspondence-based statistical shape modeling (SSM) facilitates population-level morphometrics but requires a cumbersome, potentially bias-inducing construction pipeline. Traditional construction pipelines require manual and computationally expensive steps, hindering their widespread use. Furthermore, such methods utilize templates or assumptions (e.g., linearity) that can bias or limit the expressivity of the variation captured by the constructed SSM. Recent advancements in deep learning have streamlined this process in inference by providing SSM prediction directly from unsegmented medical images. However, the proposed approaches are fully supervised and require utilizing a traditional SSM construction pipeline to create training data, thus inheriting the associated burdens and limitations. To address these challenges, we introduce a weakly supervised deep learning approach to predict SSM from images using point cloud supervision. Specifically, we propose reducing the supervision associated with the state-of-the-art fully Bayesian variational information bottleneck DeepSSM (BVIB-DeepSSM) model. BVIB-DeepSSM is an effective, principled framework for predicting probabilistic anatomical shapes from images with quantification of both aleatoric and epistemic uncertainties. Whereas the original BVIB-DeepSSM method requires strong supervision in the form of ground truth correspondence points, the proposed approach utilizes weak supervision via point cloud surface representations, which are more readily obtainable. Furthermore, the proposed approach learns correspondence in a completely data-driven manner without prior assumptions about the expected variability in shape cohort. Our experiments demonstrate that this approach yields similar accuracy and uncertainty estimation to the fully supervised scenario while substantially enhancing the feasibility of model training for SSM construction.

## 1 Introduction

Statistical shape modeling (SSM) has emerged as a useful tool in medical imaging and computational anatomy, offering valuable insights into the variability of anatomical structures, such as organs or bones, across a given population.

SSM provides a population-level statistical representation of morphology, enabling wide-ranging applications in clinical research, including disease diagnosis [10], treatment planning [22], surgical simulation [21], and outcome prediction [7]. In SSM, shapes are either represented explicitly via landmark or correspondence points, or implicitly via deformation fields (coordinate transformations in relation to a predefined or learnable atlas) [28,15]. The point distribution model (PDM) is a widely adopted explicit shape representation consisting of dense sets of correspondence points defined on the surface of the anatomical shapes in semantically consistent locations across the population. Traditionally, PDMs were automatically defined on preprocessed shape cohorts (segmented from medical images) via pairwise mapping to a predefined or learned atlas/template (e.g., [34]) or via groupwise optimization (e.g., [13]). Such SSM construction pipelines require time-consuming and expert-driven steps such as segmentation, shape registration, and optimization parameter tuning or atlas construction. Furthermore, each time a new shape is added, the pipeline must be rerun as optimization is performed across the entire cohort simultaneously.

Deep learning approaches, such as DeepSSM [12,11], offer an alternative to traditional pipelines by leveraging trained neural networks to directly infer PDMs from unsegmented volumetric images with minimal preprocessing. In inference, this alleviated the need for segmentation and reoptimization given a new scan. However, integrating deep learning-based solutions into clinical practice necessitates understanding the uncertainty associated with model predictions. Therefore, Bayesian deep learning frameworks have been proposed to provide probabilistic PDM predictions capable of quantifying the two primary forms of uncertainty: aleatoric (data-dependent) and epistemic (model-dependent) [1,2,37,6]. A notable approach is BVIB-DeepSSM [6], a probabilistic formulation of DeepSSM [11] that utilizes a fully Bayesian extension of the variational information bottleneck (VIB) framework [8]. BVIB-DeepSSM provides PDM prediction from unsegmented images with estimates of aleatoric and epistemic uncertainty that correlate with prediction error, ensuring reliable prediction without compromising accuracy.

Even though deep learning approaches mitigate the overhead associated with SSM construction during inference, they still depend on traditional SSM techniques to construct image/PDM pairs to supervise network training. This reliance not only slows down the training preparation process but also means that the network inherits any limiting assumptions made during the construction of training PDMs. Such biases or assumptions can arise from various sources, such as atlas selection in pairwise surface matching approaches or in the definition of optimization objectives. For instance, the current state-of-the-art (SOTA) groupwise optimization PDM construction method, known as particle-based shape modeling (PSM) [13,20], imposes a linearity assumption. This assumption restricts the ability of PSM to accurately represent complex, nonlinear shape variations. Training networks on PSM-constructed PDMs could similarly bias network predictions.

We propose leveraging weak supervision from point clouds in BVIB-DeepSSM training to overcome these limitations. Point cloud shape representations consist of sets of unordered, nonuniform points that sample the surface of the shape. Recently, there has been growing interest in learning SSM from point clouds due to their ease of acquisition compared to the complete, noise-free surface representations (such as meshes or binary volumes) required by traditional SSM construction methods [5,3]. This work proposes training BVIB-DeepSSM using image/point cloud pairs instead of image/PDM pairs. This approach significantly reduces the required supervision and enables training the model on readily available segmentation datasets. Our contributions are summarized as follows:

- We provide a framework to predict SSM directly from images with reduced supervision by utilizing point cloud shape representations in training rather than ground truth PDMs.
- We introduce formulations of the VIB and fully Bayesian VIB objectives that utilize permutation-invariant Chamfer distance.
- We provide comprehensive experiments that demonstrate that the proposed approach improves the feasibility of predicting SSM from images without sacrificing accuracy or uncertainty calibration.

## 2 Related Work

Traditional PDM construction methods utilize metrics such as entropy [14] or minimum description length [16], or employ parametric representations [32,35,31]. PSM [13] represents the SOTA optimization-based technique for group-wise SSM construction [20]. However, PSM assumes linear correlations, leading to a bias in the captured population variation.

DeepSSM [12,11] was the pioneering deep learning approach to predict PDMs directly from raw, unsegmented images. DeepSSM utilizes PDM supervision, where training labels are constructed via the full PSM pipeline (including segmentation, preprocessing and alignment, and PDM optimization). Uncertain-DeepSSM [1] adapted the DeepSSM network to be Bayesian, providing aleatoric and epistemic uncertainties. DeepSSM, Uncertain-DeepSSM, and other formulations [37] rely on a supervised low-dimensional encoding (i.e., shape descriptors), precomputed using principal component analysis (PCA). PCA supervision enforces a linear relationship between the latent and the output spaces and restricts the learning task to strictly SSM prediction. Additionally, PCA does not scale in the case of large sets of high-dimensional shape data. In contrast, VIB-DeepSSM [2] introduced a variational information bottleneck (VIB) architecture [8] to learn a low-dimensional latent encoding tailored to the PDM estimation task, leading to improved generalization and more accurate estimation of aleatoric uncertainty. However, the VIB framework is only half Bayesian [9]; thus, VIB-DeepSSM lacks the capability to quantify epistemic uncertainty. BVIB-DeepSSM [6] extended the VIB-DeepSSM framework to be fully Bayesian, enabling the prediction of probabilistic shapes from images with quantification for both forms of uncertainty. This SOTA model is the basis of the proposed approach.

Recent work has explored unsupervised estimation of SSM from various shape representations [5,3,24,23]. One study demonstrated that networks designed for point cloud competition perform reasonably well at the task of anatomical PDM generation [5]. These networks typically have an encoder-decoder architecture with a bottleneck [17]. The decoder provides a continuous mapping from the learned latent space to output space, resulting in consistently ordered output point clouds, providing correspondence as a by-product. Mesh2SSM [24] explicitly predicts PDMs in an unsupervised manner from mesh shape representations, utilizing complete surface information. Point2SSM [3] is a self-supervised technique proposed to predict anatomical SSM from point clouds. By employing Chamfer distance reconstruction loss, Point2SSM encourages the predicted PDMs to accurately sample the entire point clouds. Recently, SCoRP [23] proposed leveraging a shape prior learned from surface meshes to predict PDMs from unsegmented images within a student/teacher framework. While closely related to the proposed method, SCoRP lacks uncertainty quantification and necessitates complete mesh surface representations. The proposed method requires only image/point cloud pairs to supervise network training, providing PDM prediction and granular uncertainty estimates in inference.

### 3 Background

#### 3.1 Notation

Let  $\mathbf{X}$ ,  $\mathbf{Y}$ , and  $\mathbf{Z}$  denote random variables and let  $\mathbf{x}$ ,  $\mathbf{y}$ , and  $\mathbf{z}$  denote realizations of those respective random variables. Given an unsegmented volumetric image of an anatomy, denoted  $\mathbf{x} \in \mathbb{R}^{H \times W \times D}$ , the goal is to predict a PDM denoted  $\hat{\mathbf{y}} \in \mathbb{R}^{3M}$ . Each PDM is a set  $M$  *ordered correspondence points*, where a 3D vector of coordinates defines each correspondence point. Training the network requires a set of paired data, denoted  $\mathcal{D} = \{\mathcal{X}, \mathcal{Y}\}$ . Here  $\mathcal{X} = \{\mathbf{x}_n\}_{n=1}^N$  is a set of  $N$  unsegmented volumetric images. In previously proposed fully supervised settings,  $\mathcal{Y} = \{\mathbf{y}_n\}_{n=1}^N$  where  $\mathbf{y}_n \in \mathbb{R}^{3M}$  denotes a ground truth PDM, constructed via a traditional pipeline, comprised of  $M$  *ordered correspondence points*. In the proposed weakly supervised setting,  $\mathbf{y}_n$  denotes a point cloud shape representation, meaning an *unordered set* of points on the surface of the shape  $n$ . The VIB framework utilizes a learned stochastic latent encoding:  $\mathcal{Z} = \{\mathbf{z}_n\}_{n=1}^N$ , where  $\mathbf{z}_n \in \mathbb{R}^L$  and  $L \ll 3M$ .

#### 3.2 Variational Information Bottleneck

In information bottleneck (IB) theory, a stochastic encoding  $\mathbf{Z}$  is learned to capture the minimal sufficient statistics required of input  $\mathbf{X}$  to predict the output  $\mathbf{Y}$  [36]. The encoding  $\mathbf{Z}$  and model parameters  $\Theta$  are estimated by maximizing the IB objective:

$$\operatorname{argmax}_{\Theta} I(\mathbf{Y}, \mathbf{Z}; \Theta) - \beta I(\mathbf{X}, \mathbf{Z}; \Theta) \quad (1)$$

where  $I$  denotes mutual information and  $\beta$  is a Lagrangian multiplier.

The first term in Eq. 1 encourages  $\mathbf{Z}$  to be maximally expressive of  $\mathbf{Y}$ , encouraging predictive accuracy. The second term encourages  $\mathbf{Z}$  to be maximally compressive of  $\mathbf{X}$ , affecting the model complexity.

In the deep variational information bottleneck (VIB) [8] approach, the IB model is parameterized via a neural network with weights  $\Theta = \{\phi, \theta\}$ , and a latent distribution is learned by minimizing the IB objective (Eq. 1). Direct calculation of mutual information is intractable in this context, and thus VIB employs variational inference to derive a theoretical lower bound on the IB objective:

$$\mathcal{L}_{VIB} = \mathbb{E}_{\hat{\mathbf{Z}} \sim q(\mathbf{Z}|\mathbf{X}, \phi)} \left[ -\log p(\mathbf{Y}|\hat{\mathbf{Z}}, \theta) \right] + \beta \text{KL} [q(\mathbf{Z}|\mathbf{X}, \phi) \| p(\mathbf{Z})] \quad (2)$$

The first term, the negative log-likelihood (NLL), encourages  $\mathbf{Z}$  to be predictive of  $\mathbf{Y}$ . The second term, the Kullback–Leibler (KL) divergence, encourages  $\mathbf{Z}$  to be compressive of  $\mathbf{X}$ . The  $\beta$  hyper-parameter controls the tradeoff.

### 3.3 BVIB-DeepSSM

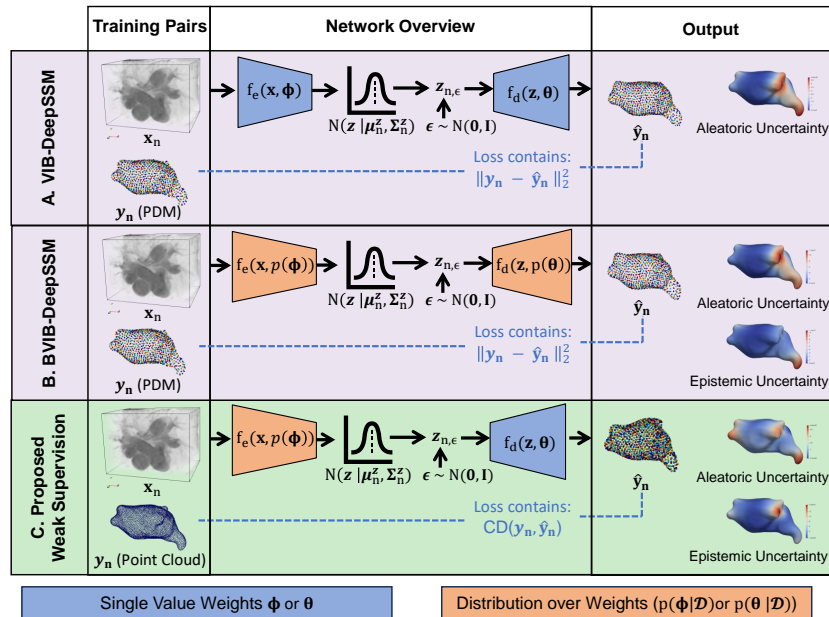
VIB-DeepSSM [2] employs the VIB [8] approach to learn the latent encoding in the context of the task: predicting PDM  $\hat{\mathbf{y}}$  from image  $\mathbf{x}$ . The VIB-DeepSSM architecture (Fig. 1.A) is comprised of an encoder and decoder. The encoder,  $f_e$ , comprised of 3D convolutional and densely connected layers parameterized by  $\phi$ , maps the input image to a Gaussian latent distribution:  $\mathcal{N}(\mathbf{z}|\mu^z, \Sigma^z)$ . Posterior samples  $\mathbf{z}_\epsilon$  are acquired from this predicted latent distribution using the reparameterization trick to enable gradient calculation. The decoder,  $f_d$  parameterized by  $\theta$ , maps the latent encoding to the predicted output  $\hat{\mathbf{y}}$ .

VIB-DeepSSM allows for capturing aleatoric uncertainty as the variance of the  $p(\mathbf{y}|\mathbf{z})$  distribution. This variance is computed by sampling multiple latent encodings from  $\mathcal{N}(\mathbf{z}|\mu^z, \Sigma^z)$  and passing them through the decoder to get a sampled distribution of predictions. A Gaussian distribution is estimated from these samples denoted:  $\mathcal{N}(\hat{\mathbf{y}}|\mu^y, \Sigma^y)$ . The estimated  $\Sigma^y$  captures the aleatoric or data-dependent uncertainty. However, this approach does not quantify epistemic uncertainty because VIB is only half-Bayesian [9].

BVIB-DeepSSM [6] derived the fully Bayesian VIB formulation by applying an additional PAC-Bound with respect to the network parameters. In the VIB-DeepSSM model, parameters  $\Theta = \{\phi, \theta\}$  are fit via maximum likelihood estimation. BVIB-DeepSSM utilizes variational inference to approximate the posterior  $p(\Theta|\mathcal{D})$ . The BVIB-DeepSSM objective results in two intractable posteriors  $p(\mathbf{Z}|\mathbf{X}, \phi)$  and  $p(\Theta|\mathbf{X}, \mathbf{Y})$ , the former is approximated via  $q(\mathbf{Z}|\mathbf{X}, \phi)$  as in Eq. 2 and the latter is approximated by  $q(\Theta)$ . The two KL divergence terms are minimized via a joint evidence lower bound, resulting in the objective:

$$\mathcal{L}_{BVIB} = \mathbb{E}_{\tilde{\Theta}} \left[ \mathbb{E}_{\hat{\mathbf{Z}}} \left[ -\log p(\mathbf{Y}|\hat{\mathbf{Z}}, \tilde{\theta}) \right] + \beta \text{KL} \left[ q(\mathbf{Z}|\mathbf{X}, \tilde{\phi}) \| p(\mathbf{Z}) \right] \right] + \text{KL} [q(\Theta) \| p(\Theta)] \quad (3)$$

In BVIB-DeepSSM (Fig. 1.B), epistemic uncertainty is captured as the variance in predictions made with various network weights sampled from the learned distribution  $p(\theta|\mathcal{D})$ . Many approaches have been proposed for the computationally challenging task of estimating a distribution over model parameters for epistemic uncertainty quantification. Among these approaches, concrete dropout and ensembling have been shown to be most effective [4] and are utilized in the BVIB-DeepSSM formulation [6]. Concrete dropout employs Monte Carlo dropout sampling as a practical approach for approximate variational inference [18]. This approach utilizes a continuous relaxation of the Bernoulli distribution (i.e., concrete distribution) to parameterize the learned distribution over weights. Epistemic uncertainty is captured by the spread of predictions resulting from inference with various sampled dropout masks. Concrete dropout automatically optimizes the dropout probabilities at each layer alongside the network weights. Batch ensemble [38] compromises between a single network and a full deep ensemble, achieving a balance between performance, computation time, and memory usage. BVIB-DeepSSM additionally proposed combining concrete dropout and batch ensemble to acquire a multimodal approximate posterior on weights for increased flexibility and expressiveness [6].



**Fig. 1.** Overview of the differences between VIB-DeepSSM [2], BVIB-DeepSSM [6], and the proposed weakly supervised variant of BVIB-DeepSSM with point cloud supervision.

## 4 Methods

We propose using point cloud shape representations to weakly supervise BVIB-DeepSSM. Recent work has shown that bottleneck network architectures with fixed decoders supervised by point cloud-based loss can learn correspondence [5]. In such networks, the bottleneck captures a population-specific shape prior. Directly decoding the latent shape feature representation results in a consistent ordering of the output point clouds across samples, providing PDMs. We propose to leverage this effect in BVIB-DeepSSM, allowing for the replacement of ground truth PDMs with unordered point clouds in the training data. This advancement requires updating the BVIB-DeepSSM formulation in two crucial ways: first, the objective must be altered for point cloud supervision, and second, consideration regarding the epistemic uncertainty quantification approach must be made.

### 4.1 Proposed Weakly Supervised Loss

Reducing the supervision requires updating the first term in the VIB and BVIB-DeepSSM objectives (Eq. 3 and 3), the NLL term. In the PDM-supervised setting, the NLL term is expressed as:

$$-\log p(\mathbf{y}|\mu^{\mathbf{y}}, \Sigma^{\mathbf{y}}) = \frac{\|\mu^{\mathbf{y}} - \mathbf{y}\|_2}{2\Sigma^{\mathbf{y}}} + \frac{1}{2} \log \Sigma^{\mathbf{y}} \quad (4)$$

where  $\mu^{\mathbf{y}}$  and  $\Sigma^{\mathbf{y}}$  are the mean and variance of the predicted distribution estimated using various posterior samples  $\hat{\Theta}$  and  $\hat{\mathbf{z}}$ .

Replacing  $\mathbf{y}$  (an ordered PDM) with unordered point clouds requires replacing the L2 norm in Eq. 4 with a permutation invariant distance metric. Chamfer distance is most commonly used for this purpose. The Chamfer distance from point cloud  $\mathbb{A}$  to point cloud  $\mathbb{B}$  is defined as:

$$\text{CD}(\mathbb{A} \rightarrow \mathbb{B}) = \frac{1}{|\mathbb{A}|} \sum_{\mathbf{a} \in \mathbb{A}} \min_{\mathbf{b} \in \mathbb{B}} \|\mathbf{a} - \mathbf{b}\|_2^2 \quad (5)$$

Typically, bidirectional Chamfer distance is used:

$$\text{CD}(\mathbb{A}, \mathbb{B}) = \text{CD}(\mathbb{A} \rightarrow \mathbb{B}) + \text{CD}(\mathbb{B} \rightarrow \mathbb{A}) \quad (6)$$

Note the number of points in  $\mathbb{A}$ , denoted  $|\mathbb{A}|$ , is not required to match  $|\mathbb{B}|$ .

We propose utilizing the single directional distance  $\text{CD}(\mu^{\mathbf{y}} \rightarrow \mathbf{y})$  as a replacement for the L2 norm in Eq. 4, as this is a commensurate metric that is calculated point-wise, in a permutation-invariant manner. The resulting updated NLL term is expressed as:

$$-\log p(\mathbf{y}|\mu^{\mathbf{y}}, \Sigma^{\mathbf{y}}) \approx \frac{\text{CD}(\mu^{\mathbf{y}} \rightarrow \mathbf{y})}{2\Sigma^{\mathbf{y}}} + \frac{1}{2} \log \Sigma^{\mathbf{y}} \quad (7)$$

This update enables permutation invariant point-wise error estimation. However, the single-directional CD does not ensure that the predicted PDM will sample

the entire surface well. For instance, if all points in  $\mu^{\mathcal{Y}}$  converge to a single point in  $\mathbf{y}$ ,  $\text{CD}(\mu^{\mathcal{Y}} \rightarrow \mathbf{y})$  would be minimized. To prevent this behavior, we include  $\text{CD}(\mathbf{y} \rightarrow \mu^{\mathcal{Y}})$  as a regularization term. This term encourages the predicted points to be well-spread across the surface so that each point in  $\mathbf{y}$  has a close neighbor in  $\mu^{\mathcal{Y}}$ . The resulting updated VIB-DeepSSM objective is expressed as: Thus, the proposed BVIB-DeepSSM objective is expressed as:

$$\mathcal{L}_{\text{Proposed VIB}} = \mathbb{E}_{\hat{\mathbf{Z}}} \left[ -\log p(\mathbf{Y}|\hat{\mathbf{Z}}, \boldsymbol{\theta}) \right] + \beta \text{KL} [q(\mathbf{Z}|\mathbf{X}, \phi) \| p(\mathbf{Z})] + \alpha \text{CD}(\mathbf{Y} \rightarrow \hat{\mathbf{Y}}) \quad (8)$$

where  $\log p(\mathbf{Y}|\hat{\mathbf{Z}}, \boldsymbol{\theta})$  is computed via Eq. 7 and  $\text{CD}(\mathbf{Y} \rightarrow \hat{\mathbf{Y}})$  is the regularization term computed via Eq. 5, weighted by hyperparameter  $\alpha$ .

## 4.2 Weakly Supervised Epistemic Uncertainty Quantification

In addition to updating the learning objective, the reduction in supervision requires adapting the approach to epistemic uncertainty quantification. While ensembling proved to be an effective frequentist approximation for estimating a distribution over model parameters in the original BVIB-DeepSSM formulation, it is not appropriate in the weakly supervised setting. This is because point cloud supervision does not enforce one particular point ordering in correspondence prediction as PDM supervision does. Rather, network-specific correspondence is induced by two factors: the Chamfer distance reconstruction loss and the consistent, continuous mapping from the latent space to the output space provided by the decoder. Thus, while a given network provides correspondence across predictions, there is no mechanism to enforce correspondence consistency across different networks or ensemble members. Each member would learn a unique output point ordering, rendering the ensemble averaging effect meaningless. Thus, in the weakly supervised context, a true Bayesian approximation method must be used to learn a distribution over weights within a single network.

Additionally, introducing stochasticity to the decoder would be detrimental to PDM prediction, as correspondence is induced by the established continuous mapping from the latent to output space. Thus, we propose adapting the concrete dropout-based BVIB-DeepSSM model to estimate epistemic uncertainty from the encoder alone and utilize a fully deterministic decoder. Here, predictive distributions are acquired by decoding various  $\mathbf{Z}$  samples with various encoder dropout masks.

Thus, the proposed BVIB-DeepSSM objective is expressed as:

$$\begin{aligned} \mathcal{L}_{\text{Proposed}} = \mathbb{E}_{\tilde{\phi}} \left[ \mathbb{E}_{\hat{\mathbf{Z}}} \left[ -\log p(\mathbf{Y}|\hat{\mathbf{Z}}, \boldsymbol{\theta}) \right] + \beta \text{KL} \left[ q(\mathbf{Z}|\mathbf{X}, \tilde{\phi}) \| p(\mathbf{Z}) \right] \right] \\ + \text{KL} [q(\phi) \| p(\phi)] + \alpha \text{CD}(\mathbf{Y} \rightarrow \hat{\mathbf{Y}}) \end{aligned} \quad (9)$$

where, as in Eq. 8,  $\log p(\mathbf{Y}|\hat{\mathbf{Z}}, \boldsymbol{\theta})$  is computed via Eq. 7.

An overview of VIB-DeepSSM, BVIB-DeepSSM, and the proposed weakly supervised BVIB-DeepSSM approach are provided in Fig. 1.



## 5 Experiments

### 5.1 Datasets

We utilize two challenging datasets to evaluate the proposed method: the left atrium and the liver. The left atrium dataset includes 1,096 shapes derived from cardiac late gadolinium enhancement MRI images of different atrial fibrillation patients. The images were manually segmented at the University of Utah Division of Cardiovascular Medicine with spatial resolution  $0.65 \times 0.65 \times 2.5 \text{ mm}^3$ , and the endocardium wall was used to cut off pulmonary veins. This dataset includes substantial morphological diversity in overall size, the size of the left atrium appendage, and the quantity and arrangement of pulmonary veins. Following BVIB-DeepSSM, we hold out outlier cases in the test set, selected via thresholding on an outlier degree computed on images and meshes [29]. The resulting test set includes 40 shape outliers, 78 image outliers, and 92 randomly selected inlier test samples. The liver dataset contains 834 CT scans and corresponding quality-controlled segmentations from the open-source AbdomenCT-1K dataset [27]. These images vary significantly in intensity, quality, and resolution, providing a challenging test case. We randomly split the liver data 80%/10%/10% into training, validation, and test sets.

### 5.2 Experimental Setup

We compare the proposed weakly supervised adaptations of VIB-DeepSSM and BVIB-DeepSSM with the original formulations. We utilize the PSM construction method implemented in the ShapeWorks software suite [13] to create PDMs for training fully supervised methods. Additionally, we utilize ShapeWorks to process images and segmentations, including cropping around the region of interest and downsampling to manage memory usage. We generate surface meshes with 5000 vertices from the segmentations. The vertices serve as point clouds for the proposed weak supervision. We employ image augmentation in training all models in the form of additive Gaussian noise with random variance between 0 and 1% of the full signal. Following the BVIB-DeepSSM strategy, burn-in is used to convert the loss from deterministic (L2 or CD) to probabilistic (Eqs 3 and 9) [6]. This burn-in counteracts the accuracy reduction that occurs when NLL-based loss is used with a gradient-based optimizer [33]. The concrete dropout implementation of BVIB-DeepSSM is used with initial dropout probabilities of 0.1. All models were trained until the validation error (either L2 or CD, depending on supervision) had not decreased in 50 epochs. The training was done on Tesla V100 GPU with Xavier initialization [19], Adam optimization [25]. Full model parameters and training and evaluation code are provided at <https://github.com/jadie1/Weakly-Supervised-BVIB-DeepSSM/>.

### 5.3 Evaluation Metrics

There are three factors to consider when evaluating probabilistic PDM prediction accuracy. The first is surface sampling accuracy, which assesses how well the

points are constrained to capture the complete shape surfaces. The second is the assessment of how well the population-level statistics are captured through predicted correspondences. The third is the calibration of the uncertainty estimates. This section describes the metrics used to assess these three factors.

**Surface Sampling:** A small **Chamfer distance** between the point cloud and predicted PDM,  $CD(\hat{\mathbf{y}}, \mathbf{y})$  (Eq. (6)), indicates the output points accurately capture the complete shape. **Point-to-surface distance (P2S)** assesses how well points are constrained to the surface. P2S quantifies the distance of the predicted points to a complete ground truth surface shape representation (i.e., mesh).

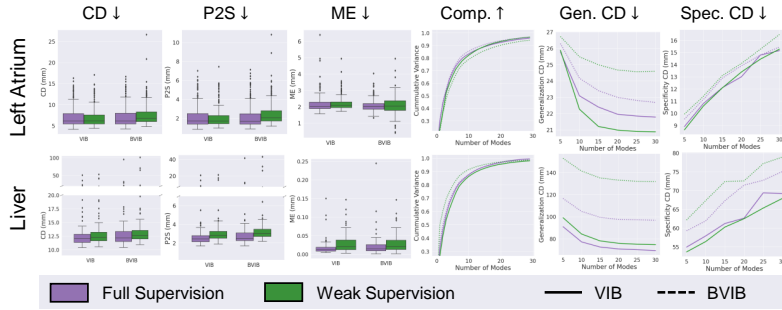
**Correspondence/SSM Metrics:** Principal component analysis (PCA) is used in SSM analysis to understand the modes of variation in the population and to evaluate how effectively population-level statistics are captured [30]. Three key metrics help measure this: compactness, generalization, and specificity. **Compactness (Comp.):** Good correspondence leads to a more compact SSM, meaning the training data distribution can be represented using a minimal number of parameters. Strong correspondence allows a larger proportion of explained population variance to be captured with fewer PCA modes. A larger area under the cumulative variance plot indicates better correspondence. **Generalization (Gen. CD):** A precise SSM should generalize effectively from training subjects to new, unseen subjects. The generalization metric measures the Chamfer distance (CD) between estimated correspondences from test point clouds and their reconstructions from training SSM-based PCA embeddings using varying numbers of components. A smaller CD indicates better generalization. **Specificity (Spec. CD):** Specificity assesses whether the predicted SSM produces valid instances of the shape class. It is calculated as the average Chamfer distance between training examples and generated samples from the training SSM-based PCA embeddings using different numbers of components. A smaller CD suggests the SSM is more specific. Recent work also utilizes mapping error (**ME**) to estimate correspondence accuracy [3]. ME quantifies how consistent output point neighborhoods are across the population [26]. A lower ME indicates consistent neighborhoods, implying better correspondence.

**Uncertainty Calibration** We expect well-calibrated uncertainty estimates to correlate with the error. Thus, a higher Pearson  $r$  coefficient between predicted uncertainty and P2F error suggests better calibration. Furthermore, accurate uncertainty estimation is useful in out-of-distribution (OOD) detection.

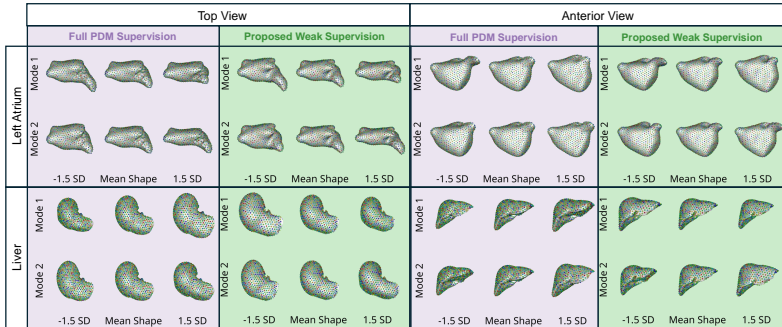
## 5.4 Results

Fig. 2 provides the surface sampling and correspondence evaluation metrics across both test sets. **The proposed weakly supervised models provide**

similar accuracy across all metrics while significantly reducing the supervision requirement. Hence, we are not sacrificing accuracy, but we democratize building networks that provide PDMs directly from unsegmented images by requiring only the level of supervision needed to train segmentation networks.



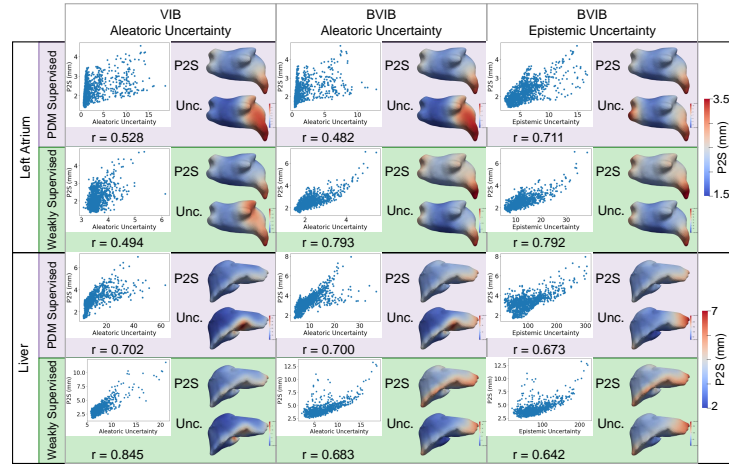
**Fig. 2.** Results of VIB-DeepSSM and C-BVIB-DeepSSM with PDM-supervision and the proposed point cloud supervision on the left atrium and liver dataset. Box plots show the distribution of errors across the test set. The SSM metric plots show the values with various numbers of PCA modes, where the line type depicts the type of DeepSSM model. Lower values are better for all metrics with the exception of compactness.



**Fig. 3.** Modes of variation resulting from the predicted SSM on the test set with the BVIB-DeepSSM models from the top and anterior view. The mean shape is shown with the primary and secondary PCA modes of variation at  $\pm 1.5$  standard deviations (SD). Correspondence points are displayed over meshes constructed from the points.

Fig. 3 displays the modes of variation resulting from the predicted SSM on the test sets for the BVIB-DeepSSM model. The mean shape and primary and secondary modes of variation resulting from PDM supervision and the proposed PC supervision are very similar. In the left atrium case, the primary mode captures the length of the left atrium appendage, and the secondary mode captures

the volume or sphericity. The primary and secondary modes of variation in the liver dataset capture the size and curvature. The results demonstrate that the proposed weak supervision does not lead to less accurate PDM prediction. The predicted points sample the true surface to the same degree and offer a similar correspondence accuracy despite the absence of ground truth correspondence supervision. Additionally, the similarity in the captured modes of variation suggests that PDM predictions made with weak supervision could be equally useful in downstream clinical tasks as the predictions made with full PDM supervision.

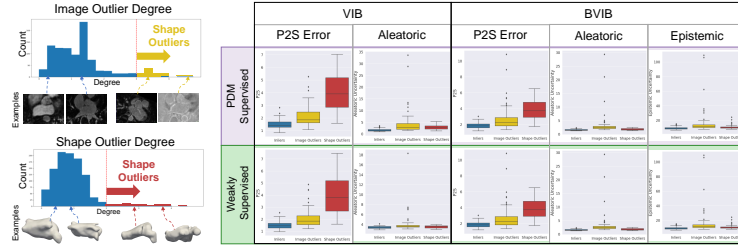


**Fig. 4.** Uncertainty calibration results. Scatter plots and corresponding Pearson R correlation coefficients demonstrate the point-wise correlation between the estimated uncertainty and P2S error across the test sets. The average P2S error and uncertainty values are also shown via heatmaps on a representative mesh, illustrating spatial correlation.

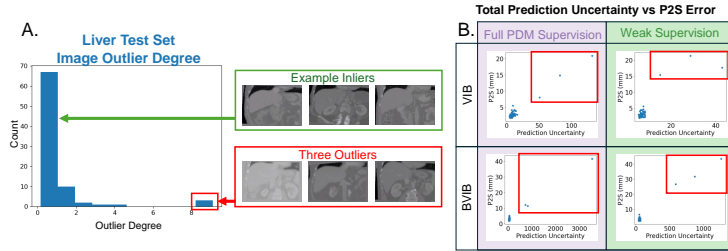
Fig. 4 displays the point-wise correlation between the predicted uncertainty values and P2S distance error across the test set. We would expect the uncertainty estimates to be higher for predicted correspondence points that are further from the true shape surface. The  $r$  correlation coefficients and resulting average uncertainty heatmaps are very similar, indicating that reducing the supervision does not significantly impact the uncertainty calibration. These results demonstrate the effectiveness of estimating the NLL term with CD rather than L2 Euclidean distance (Eq. 7). The spatial correlation between the P2S error and uncertainty heatmaps demonstrates the utility of these probabilistic frameworks in aiding in assessing prediction reliability.

Uncertainty quantification is also useful in detecting out-of-distribution (OOD) samples. The left atrium dataset is comprised of three subsets: image outliers, shape outliers, and randomly selected inlier examples. This partitioning was per-

formed by thresholding on a precomputed outlier degree [29] as shown in Fig. 5. The error and uncertainty estimation distributions across subsets are shown in Fig. 5. The predicted uncertainty is slightly higher for the outlier test sets, especially for the extreme image outliers. The full and weakly supervised models provide similar patterns in error and uncertainty across the left atrium test sets.



**Fig. 5.** Left atrium outlier test sets. The histogram plots the distribution of image and shape outlier degrees with example image slices and meshes. Box plots show the distribution of P2S error and uncertainty across the three test sets: inliers (blue), image outliers (yellow), and shape outliers (red).



**Fig. 6.** Liver outlier detection results. The histogram plots the distribution of image outlier degrees across the liver test set, highlighting the three outliers. Slices of inlier and outlier images are shown. The total prediction uncertainty and P2S error (averaged across each shape) are plotted across for each model. The three outlier cases have high uncertainty and are highlighted in red boxes.

Fig. 6 shows the distribution of image outlier degrees across the randomly selected liver test set. Three outlier cases are identified in this histogram. These three cases are also clearly identifiable in the P2S error vs prediction uncertainty scatter plots in Fig. 6. Here, prediction uncertainty is the total aleatoric for VIB models and the sum of the total aleatoric and epistemic estimates for BVIB. The outlier cases are clearly identifiable, given the prediction uncertainty resulting from both full and weak supervision, suggesting the proposed weakly supervised approach is not detrimental to OOD detection.

## 6 Conclusion

We proposed an alternative training approach to BVIB-DeepSSM with reduced supervision. The proposed framework matches PDM-supervision accuracy while significantly streamlining the training pipeline. In future work, the point cloud shape representations could also be leveraged to learn a more expressive prior  $p(\mathbf{z})$ . In [39], it is proven that learning the variational autoencoder (VAE) latent prior is necessary for reaching the extremum of the VAE objective. This proof can be directly applied to show learning  $p(\mathbf{z})$  is necessary for reaching the extremum of the VIB objective. Future work could explore utilizing a point cloud autoencoder to learn  $p(\mathbf{z})$  in a shape-informed manner. Overall, this work improves the feasibility of SSM construction from images, making SSM more accessible as a tool for clinical research.

## 7 Acknowledgements

This work was supported by the National Institutes of Health under grant numbers NIBIB-U24EB029011, NIAMS-R01AR076120, NHLBI-R01HL135568, and NIBIB-R01EB016701. We thank the University of Utah Division of Cardiovascular Medicine for providing left atrium MRI scans and segmentations from the Atrial Fibrillation projects and the ShapeWorks team.

## References

1. Adams, J., Bhalodia, R., Elhabian, S.: Uncertain-deepssm: From images to probabilistic shape models. In: Shape in Medical Imaging: International Workshop, ShapeMI 2020, Held in Conjunction with MICCAI 2020, Lima, Peru, October 4, 2020, Proceedings. pp. 57–72. Springer (2020)
2. Adams, J., Elhabian, S.: From images to probabilistic anatomical shapes: A deep variational bottleneck approach. In: International Conference on Medical Image Computing and Computer-Assisted Intervention. pp. 474–484. Springer (2022)
3. Adams, J., Elhabian, S.: Point2ssm: Learning morphological variations of anatomies from point cloud. arXiv preprint arXiv:2305.14486 (2023)
4. Adams, J., Elhabian, S.Y.: Benchmarking scalable epistemic uncertainty quantification in organ segmentation. In: International Workshop on Uncertainty for Safe Utilization of Machine Learning in Medical Imaging. pp. 53–63. Springer (2023)
5. Adams, J., Elhabian, S.Y.: Can point cloud networks learn statistical shape models of anatomies? In: International Conference on Medical Image Computing and Computer-Assisted Intervention. pp. 486–496. Springer (2023)
6. Adams, J., Elhabian, S.Y.: Fully bayesian vib-deepssm. In: International Conference on Medical Image Computing and Computer-Assisted Intervention. pp. 346–356. Springer (2023)
7. Aldieri, A., Bhattacharya, P., Paggiosi, M., Eastell, R., Audenino, A.L., Bignardi, C., Morbiducci, U., Terzini, M.: Improving the hip fracture risk prediction with a statistical shape-and-intensity model of the proximal femur. *Annals of biomedical engineering* **50**(2), 211–221 (2022)

8. Alemi, A.A., Fischer, I., Dillon, J.V., Murphy, K.: Deep variational information bottleneck. arXiv preprint arXiv:1612.00410 (2016)
9. Alemi, A.A., Morningstar, W.R., Poole, B., Fischer, I., Dillon, J.V.: Vib is half bayes. In: Third Symposium on Advances in Approximate Bayesian Inference (2020)
10. Ambellan, F., Lamecker, H., von Tycowicz, C., Zachow, S.: Statistical shape models: understanding and mastering variation in anatomy. Springer (2019)
11. Bhalodia, R., Elhabian, S., Adams, J., Tao, W., Kavan, L., Whitaker, R.: Deepssm: A blueprint for image-to-shape deep learning models. *Medical Image Analysis* **91**, 103034 (2024)
12. Bhalodia, R., Elhabian, S.Y., Kavan, L., Whitaker, R.T.: Deepssm: A deep learning framework for statistical shape modeling from raw images. In: *Shape In Medical Imaging at MICCAI. Lecture Notes in Computer Science*, vol. 11167, pp. 244–257. Springer (2018)
13. Cates, J., Elhabian, S., Whitaker, R.: Shapeworks: Particle-based shape correspondence and visualization software. In: *Statistical Shape and Deformation Analysis*, pp. 257–298. Elsevier (2017)
14. Cates, J., Fletcher, P.T., Styner, M., Shenton, M., Whitaker, R.: Shape modeling and analysis with entropy-based particle systems. In: *IPMI*. pp. 333–345. Springer (2007)
15. Cootes, T.F., Twining, C.J., Taylor, C.J.: Diffeomorphic statistical shape models. In: *BMVC*. pp. 1–10. Citeseer (2004)
16. Davies, R.H., Twining, C.J., Cootes, T.F., Waterton, J.C., Taylor, C.J.: A minimum description length approach to statistical shape modeling. *IEEE transactions on medical imaging* **21**(5), 525–537 (2002)
17. Fei, B., Yang, W., Chen, W.M., Li, Z., Li, Y., Ma, T., Hu, X., Ma, L.: Comprehensive review of deep learning-based 3d point cloud completion processing and analysis. *IEEE Transactions on Intelligent Transportation Systems* (2022)
18. Gal, Y., Hron, J., Kendall, A.: Concrete dropout. *Advances in neural information processing systems* **30** (2017)
19. Glorot, X., Bengio, Y.: Understanding the difficulty of training deep feedforward neural networks. In: *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics. Proceedings of Machine Learning Research*, vol. 9, pp. 249–256. PMLR (13–15 May 2010)
20. Goparaju, A., Iyer, K., Bone, A., Hu, N., Henninger, H.B., Anderson, A.E., Durrleman, S., Jacxsens, M., Morris, A., Csecs, I., et al.: Benchmarking off-the-shelf statistical shape modeling tools in clinical applications. *Medical image analysis* **76**, 102271 (2022)
21. Haq, R., Schmid, J., Borgie, R., Cates, J., Audette, M.A.: Deformable multisurface segmentation of the spine for orthopedic surgery planning and simulation. *Journal of medical imaging* **7**(1), 015002–015002 (2020)
22. Hassan, M.K., Fleury, E., Shamonin, D., Fonk, L.G., Marinkovic, M., Jaarsma-Coes, M.G., Luyten, G.P., Webb, A., Beenakker, J.W., Stoel, B.: An automatic framework to create patient-specific eye models from 3d magnetic resonance images for treatment selection in patients with uveal melanoma. *Advances in Radiation Oncology* **6**(6), 100697 (2021)
23. Iyer, K., Adams, J., Elhabian, S.Y.: Scorp: Statistics-informed dense correspondence prediction directly from unsegmented medical images. arXiv preprint arXiv:2404.17967 (2024)

24. Iyer, K., Elhabian, S.Y.: Mesh2ssm: From surface meshes to statistical shape models of anatomy. In: International Conference on Medical Image Computing and Computer-Assisted Intervention. pp. 615–625. Springer (2023)
25. Kingma, D., Ba, J.: Adam: A method for stochastic optimization. International Conference on Learning Representations (12 2014)
26. Lang, I., Ginzburg, D., Avidan, S., Raviv, D.: Dpc: Unsupervised deep point correspondence via cross and self construction. In: 2021 International Conference on 3D Vision (3DV). pp. 1442–1451. IEEE (2021)
27. Ma, J., Zhang, Y., Gu, S., Zhu, C., Ge, C., Zhang, Y., An, X., Wang, C., Wang, Q., Liu, X., Cao, S., Zhang, Q., Liu, S., Wang, Y., Li, Y., He, J., Yang, X.: Abdoment-1k: Is abdominal organ segmentation a solved problem? IEEE Transactions on Pattern Analysis and Machine Intelligence **44**(10), 6695–6714 (2022). <https://doi.org/10.1109/TPAMI.2021.3100536>
28. Miller, M.I., Younes, L., Trouvé, A.: Diffeomorphometry and geodesic positioning systems for human anatomy. Technology **2**(01), 36–43 (2014)
29. Moghaddam, B., Pentland, A.: Probabilistic visual learning for object representation. IEEE Transactions on pattern analysis and machine intelligence **19**(7), 696–710 (1997)
30. Munsell, B.C., Dalal, P., Wang, S.: Evaluating shape correspondence for statistical shape analysis: A benchmark study. IEEE Transactions on Pattern Analysis and Machine Intelligence **30**(11), 2023–2039 (2008)
31. Nain, D., Styner, M., Niethammer, M., Levitt, J.J., Shenton, M.E., Gerig, G., Bobick, A., Tannenbaum, A.: Statistical shape analysis of brain structures using spherical wavelets. In: 2007 4th IEEE International Symposium on Biomedical Imaging: From Nano to Macro. pp. 209–212. IEEE (2007)
32. Ovsjanikov, M., Ben-Chen, M., Solomon, J., Butscher, A., Guibas, L.: Functional maps: a flexible representation of maps between shapes. ACM Transactions on Graphics (ToG) **31**(4), 1–11 (2012)
33. Seitzer, M., Tavakoli, A., Antic, D., Martius, G.: On the pitfalls of heteroscedastic uncertainty estimation with probabilistic neural networks. In: International Conference on Learning Representations (2021)
34. Styner, M., Oguz, I., Xu, S., Brechbühler, C., Pantazis, D., Levitt, J.J., Shenton, M.E., Gerig, G.: Framework for the statistical shape analysis of brain structures using spharm-pdm. The insight journal p. 242 (2006)
35. Styner, M., Oguz, I., Xu, S., Brechbühler, C., Pantazis, D., Levitt, J.J., Shenton, M.E., Gerig, G.: Framework for the statistical shape analysis of brain structures using spharm-pdm. The insight journal p. 242 (2006)
36. Tishby, N., Pereira, F.C., Bialek, W.: The information bottleneck method. arXiv preprint physics/0004057 (2000)
37. Tóthová, K., Parisot, S., Lee, M.C., Puyol-Antón, E., Koch, L.M., King, A.P., Konukoglu, E., Pollefeys, M.: Uncertainty quantification in cnn-based surface prediction using shape priors. In: Shape in Medical Imaging: International Workshop, ShapeMI 2018, Held in Conjunction with MICCAI 2018, Granada, Spain, September 20, 2018, Proceedings. pp. 300–310. Springer (2018)
38. Wen, Y., Tran, D., Ba, J.: Batchensemble: an alternative approach to efficient ensemble and lifelong learning. In: International Conference on Learning Representations (2020)
39. Xu, H., Chen, W., Lai, J., Li, Z., Zhao, Y., Pei, D.: On the necessity and effectiveness of learning the prior of variational auto-encoder. arXiv preprint arXiv:1905.13452 (2019)