

# SCorP: Statistics-Informed Dense Correspondence Prediction Directly from Unsegmented Medical Images

Krithika Iyer<sup>1,2</sup>[0000–0003–2295–8618], Jadie Adams<sup>1,2</sup>[0000–0001–7774–5148], and Shireen Y. Elhabian<sup>1,2</sup>[0000–0002–7394–557X]

<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  
krithika.iyer@utah.edu {jadie,shireen}@sci.utah.edu

**Abstract.** Statistical shape modeling (SSM) is a powerful computational framework for quantifying and analyzing the geometric variability of anatomical structures, facilitating advancements in medical research, diagnostics, and treatment planning. Traditional methods for shape modeling from imaging data demand significant manual and computational resources. Additionally, these methods necessitate repeating the entire modeling pipeline to derive shape descriptors (e.g., surface-based point correspondences) for new data. While deep learning approaches have shown promise in streamlining the construction of SSMs on new data, they still rely on traditional techniques to supervise the training of the deep networks. Moreover, the predominant linearity assumption of traditional approaches restricts their efficacy, a limitation also inherited by deep learning models trained using optimized/established correspondences. Consequently, representing complex anatomies becomes challenging. To address these limitations, we introduce SCorP, a novel framework capable of predicting surface-based correspondences directly from unsegmented images. By leveraging the shape prior learned directly from surface meshes in an unsupervised manner, the proposed model eliminates the need for an optimized shape model for training supervision. The strong shape prior acts as a teacher and regularizes the feature learning of the student network to guide it in learning image-based features that are predictive of surface correspondences. The proposed model streamlines the training and inference phases by removing the supervision for the correspondence prediction task while alleviating the linearity assumption. Experiments on the LGE MRI left atrium dataset and Abdomen CT-1K liver datasets demonstrate that the proposed technique enhances the accuracy and robustness of image-driven SSM, providing a compelling alternative to current fully supervised methods.

**Keywords:** Statistical Shape Modeling · Representation Learning · Correspondence Models · Deep Learning

## 1 Introduction

Statistical shape modeling (SSM) is a computational approach for statistically representing anatomies in the context of a population. SSM finds diverse applications in biomedical research, from visualizing organs [8], bones [27], and tumors [22], to assisting in surgical planning [24], disease monitoring [31], and implant design [14]. Shapes can be represented *explicitly* by a set of ordered landmarks or *correspondence* points, aka point distribution models (PDMs), or implicitly using techniques such as deformation fields [13] or level sets [25]. This paper focuses on explicit shape representations (i.e., PDMs), characterized by a dense set of correspondences describing anatomically equivalent points across samples. PDM is preferred for its simplicity and efficacy in facilitating interpretable shape comparisons and statistical analyses across populations [10].

State-of-the-art SSM methods typically require a labor-intensive and computationally demanding workflow that includes manual segmentation of anatomical structures, requiring specialized expertise. Segmentation is followed by pre-processing (e.g., resampling, cropping, and shape registration) and correspondence optimization. This entire process has to be repeated at inference (i.e., for new images), hindering feasibility as an on-demand diagnostic tool in clinical settings. Deep learning models have emerged as alternatives to traditional tools. Models such as DeepSSM and TL-DeepSSM [7,6] learn to estimate correspondences from unsegmented CT/MRI images. Despite their potential, these deep learning approaches still rely on supervised losses and necessitate established PDMs from traditional methods for training. This dependence on established PDMs introduces linearity assumptions, affecting the ability of the models to represent complex anatomical structures adequately. Additionally, this burdensome training requirement inhibits the models’ scalability and generalization. Furthermore, such deep-learning models depend on shape-based generative data augmentation strategies (via principal component analysis (PCA), non-parametric kernel density estimation (KDE), or Gaussian mixture models), requiring extensive offline computation and imposing time burden.

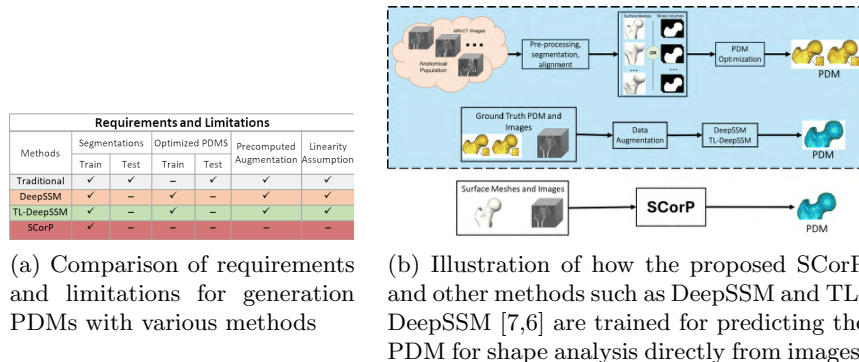


Fig. 1: Comparison of requirements and training pipelines

The newest breed of SSM deep learning models remove linearity assumptions and drop the requirement for the ground truth PDMs for training [3,17]. Despite these efforts to improve SSM methodologies, using images directly to predict the correspondences remains challenging. The inherent challenge lies in achieving high-quality shape correspondences from unsegmented images, and supervising the training of these models using an established PDM remains a bottleneck. To tackle these challenges, we propose a novel deep-learning model SCorP that is capable of predicting correspondences directly from images by leveraging shape prior built directly from the surface representation of anatomies. The shape prior can be learned from different surface representations encompassing various forms such as meshes, point clouds, and segmentations, thereby enhancing the model’s versatility and applicability.

Volumetric images (e.g., CT/MRI) may present challenges, including (a) noisy and unreliable image features like intensity and texture and (b) poorly defined anatomy boundaries, particularly in low-contrast environments. Furthermore, images depicting irregular shapes with high variability can impede the identification of invariant features. However, when specific anatomical classes are anticipated, integrating shape prior information can guide and constrain the correspondence estimation process to overcome these challenges. Our proposed model SCorP takes advantage of existing multi-view data, consisting of paired volumetric images and surface representations, through a *teacher and student* framework. In this framework, shape prior serves as the *teacher* for image-based learning. By guiding the *student* network responsible for feature extraction in the image-driven SSM task, enhancing accuracy and robustness.

Figure 1 a and b provide an overview of the requirements, limitations, and visual comparison of different SSM pipelines. Notably, our proposed method SCorP distinguishes itself by its minimal requirements (Figure 1.a), relying solely on surface representation in the form of meshes, point clouds, or binary volumes for training while avoiding adherence to the linearity assumption. Figure 1.b further provides a visual comparison of the training pipelines of various methods in contrast to SCorP. Our main contributions are:

1. We introduce **S**tatistics-informed **C**orrespondence **P**rediction (**SCorP**), a novel deep learning model designed to predict shape correspondences directly from images. By leveraging the statistics learned from surface representations as a shape prior, our model enables accurate inference of shape descriptors directly from images, bypassing the need for optimization and parameter tuning required in traditional methods.
2. We validate the accuracy of SCorP through experiments conducted on the CT (AbdomenCT-1K liver) dataset [21] and LGE MRI (left atrium) dataset. Furthermore, experiments involving varying training dataset sizes provide evidence of the model’s robustness and generalization capabilities.

## 2 Related Work

Various traditional methods for establishing correspondences have been proposed, including non-optimized landmark estimation through warping an anno-

tated reference using registration [16], parametric methods using basis functions [26], and non-parametric optimization techniques (e.g., particle-based optimization [9] and minimum description length (MDL) [12]). Non-optimized and parametric methods fail to handle complex shapes due to their fixed geometric basis or predefined template. Non-parametric optimization methods offer a more robust approach by considering the variability of the entire cohort during optimization but still rely on limiting assumptions to define optimization objective (i.e., linearity).

Deep learning models such as DeepSSM and TL-DeepSSM [7,6] provide alternatives to traditional SSM tools and are gaining traction. These models learn a functional mapping parameterized by a deep network that estimates surface correspondences from unsegmented images in a supervised manner. Several models have been proposed to enhance the performance of DeepSSM [7,6]. These modifications include incorporating multi-scale and progressive learning modules (e.g., Progressive DeepSSM [5]), introducing anatomy localization modules for raw images (e.g., LocalizedSSM [28]), and introducing uncertainty quantification (e.g., Uncertain DeepSSM [1], VIB-DeepSSM [2], BVIB-DeepSSM [4]). Despite these advancements, these models still rely on optimized PDMs for training. Other deep learning-based image-driven SSM methods have been introduced that leverage radial basis functions (RBF)-based representation to learn control points and normals for surface estimation [30]. However, these models face scalability challenges with large datasets and increased correspondences required to model complex anatomies.

Among the new breed of SSM techniques, Point2SSM [3] learns correspondences from unstructured point clouds without connectivity information that represents the surface of the anatomy. However, connectivity information can provide valuable insights when dealing with complex anatomical structures, which leads us to the models that operate on the surface meshes. Models such as FlowSSM [20] and ShapeFlow [18] employ neural networks to parameterize deformation fields on surface meshes in a low-dimensional latent space, adopting an encoder-free configuration. However, these methods necessitate re-optimization for latent representations of individual mesh samples, posing a notable challenge. Mesh2SSM [17] overcomes this issue by replacing the encoder-free setup with geodesic features and EdgeConv [29] based mesh autoencoder.

In summary, this review of methods paves the way for our proposed framework, which enhances the image-driven SSM task by directly predicting correspondences from images. Our framework incorporates a principled shape prior and eliminates the need for established PDM supervision during training.

### 3 Method

This section presents the formulation, training, and inference phases of SCoRP. Comprehensive details on network architectures and implementation specifics are provided in the Appendix. Surface meshes, point clouds, and binary volumes are all viable forms of surface representation. Without loss of generality, we primarily focus on surface meshes for notation simplicity. However, any surface



representation can be used by simply using the relevant architecture for the surface encoder.

Consider a training dataset consisting of  $N$  aligned surface meshes denoted as  $\mathcal{S} = \{S_1, S_2, \dots, S_N\}$  along with their corresponding aligned volumetric images denoted as  $\mathcal{I} = \{\mathbf{I}_1, \mathbf{I}_2, \dots, \mathbf{I}_N\}$ . Each surface mesh is denoted by  $S_j = (\mathcal{V}_j, \mathcal{E}_j)$ , where  $\mathcal{V}_j$  and  $\mathcal{E}_j$  denote the vertices and edge connectivity, respectively. The primary objective of the model is to establish a shape prior i.e., *teacher*, by learning to predict a set of  $M$  correspondence points  $\mathcal{C}_j^S = \{\mathbf{c}_{j(1)}, \mathbf{c}_{j(2)}, \dots, \mathbf{c}_{j(M)}\}$  with  $\mathbf{c}_{j(m)} \in \mathbb{R}^3$  that comprehensively describe the anatomy represented by the surface mesh  $S_j$ . Subsequently, the model leverages this shape prior to guide the feature learning of the image encoder, i.e., *student*, towards extracting image features more conducive to predicting a set of correspondence  $\mathcal{C}_j^I = \{\mathbf{c}_{j(1)}, \mathbf{c}_{j(2)}, \dots, \mathbf{c}_{j(M)}\}$  with  $\mathbf{c}_{j(m)} \in \mathbb{R}^3$  directly from the associated image  $\mathbf{I}_j$ .

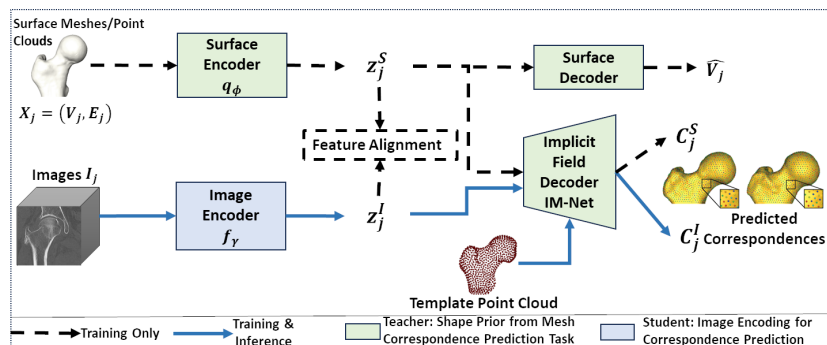


Fig. 2: **Architecture of SCorP:** Training involves three phases: (1) **Surface branch training** focuses on shape prior development using the teacher network consisting of the surface autoencoder and IM-NET decoder; (2) **Image branch embedding alignment** trains the student i.e., image encoder to predict image feature that aligns with the shape prior; (3) **Image branch prediction refinement** improves predicted correspondences from images.

### 3.1 Surface Autoencoder and Implicit Field Decoder

To learn the shape prior, i.e., *teacher*, we begin by training a surface autoencoder to learn a low-dimensional representation vector  $\mathbf{z}_j^S \in \mathbb{R}^L$  for each surface mesh  $S_j = (\mathcal{V}_j, \mathcal{E}_j)$ . We adopt state-of-the-art dynamic graph convolution that employs EdgeConv blocks [29] (akin to Mesh2SSM [17] and Point2SSM [3]) to capture permutation invariant local geometric mesh features. EdgeConv blocks compute edge features for each vertex using nearest neighbor computation. These features are then globally aggregated to produce a 1D global descriptor  $\mathbf{z}_j^S$  representing the mesh. Notably, the initial EdgeConv block utilizes geodesic distance for feature calculation on the mesh surface. In the case of point cloud data, the original architecture of EdgeConv [29] graph convolution network without the geodesic information is employed for feature extraction.

The IM-NET [11] architecture utilizes the feature vector  $\mathbf{z}_j^S$  to predict correspondences  $\mathcal{C}_j^S$  for the mesh  $S_j$ . This network uses a template point cloud and enforces correspondence relationships across samples by estimating the deformation needed for each point in the template to align with each sample based on  $\mathbf{z}_j^S$ . IM-NET transforms the template point cloud to match each sample, ensuring consistent correspondence across the dataset.

The surface autoencoder and implicit field decoder are trained jointly to minimize the two-way  $L_2$  Chamfer distance metric between the predicted correspondences  $\mathcal{C}_j^S$  and the mesh vertices  $\mathcal{V}_j$  (or point cloud coordinates when considering point clouds), and the reconstruction loss of the autoencoder between the input vertex locations  $\mathcal{V}_j$  and the reconstructed vertex locations  $\hat{\mathcal{V}}_j$ . The combined loss function  $\mathcal{L}_S$  is expressed as:

$$\mathcal{L}_S = \sum_{j=1}^N \left[ \mathcal{L}_{CD}(\mathcal{V}_j, \mathcal{C}_j^S) + \alpha \mathcal{L}_{MSE}(\mathcal{V}_j, \hat{\mathcal{V}}_j) \right] \quad (1)$$

where  $\alpha$  is the weighting factor for the vertex reconstruction term.

### 3.2 Image Encoder

The goal of the student network, i.e., the image encoder module, is to learn a compact representation  $\mathbf{z}_j^I \in \mathbb{R}^L$  for each input image  $\mathbf{I}_j$ . Like surface meshes, the latent representation  $\mathbf{z}_j^I$  will generate the correspondences. To ensure that the encoder captures meaningful representations of the underlying anatomy and is predictive of correspondences, we integrate the shape prior obtained from the teacher, i.e., the surface encoder and implicit field decoder. This integration occurs at two levels: embedding alignment and prediction refinement.

*Embedding alignment phase* aligns image features with corresponding surface features, achieved through a regression loss in both surface mesh and image embedding space. Embedding alignment teaches the image encoder to learn representations in the image domain that are semantically meaningful and coherent. Thus, the model learns to map image features to proximal mesh feature regions by minimizing the regression loss in the embedding space. The loss function for image feature alignment is denoted as:

$$\mathcal{L}_{EA} = \frac{1}{N} \sum_{j=1}^N \left[ |q_\phi(\mathbf{z}_j^S | S_j) - f_\gamma(\mathbf{z}_j^I | \mathbf{I}_j)|^2 \right] \quad (2)$$

*Prediction refinement phase* refines correspondences predicted by the image branch to match the surface mesh better. Refinement is done by minimizing the Chamfer distance between the predicted correspondences  $\mathcal{C}_j^I$  from the image  $\mathbf{I}_j$  and the mesh vertices  $\mathcal{V}_j$ . This enables the model to refine the initial predictions learned after the embedding alignment phase. The loss function for image

branch prediction refinement is denoted as:

$$\mathcal{L}_{PR} = \sum_{j=1}^N \mathcal{L}_{L_2CD}(\mathcal{V}_j, \mathcal{C}_j^I) \quad (3)$$

### 3.3 Training Strategy

SCorP’s training process involves three phases, each focusing on different aspects of the model architecture. The overall loss function guiding the training is formulated as  $\mathcal{L} = \lambda_1 \mathcal{L}_S + \lambda_2 \mathcal{L}_{EA} + \lambda_3 \mathcal{L}_{PR}$  where  $\lambda_1$  and  $\lambda_2$ , and  $\lambda_3$  are the weighting factors.

1. **Surface branch training:** We begin by training the teacher network consisting of the surface autoencoder and the implicit field decoder. This phase aims to develop a correspondence model based on surface representation, i.e., shape prior. During this phase, the loss function is defined as  $\mathcal{L} = \mathcal{L}_S$ , with  $\lambda_1 = 1$  and  $\lambda_2 = \lambda_3 = 0$ .
2. **Image Branch Embedding Alignment:** Next, we focus on training the student network, i.e., the image branch embedding, while keeping the teacher network weights unchanged. This allows the image encoder to learn a shared manifold consistent with the volumetric image and the surface representation. The loss function for this phase is  $\mathcal{L} = \mathcal{L}_{EA}$ , with  $\lambda_1 = \lambda_3 = 0$  and  $\lambda_2 = 1$ .
3. **Image Branch Prediction Refinement:** Finally, the predicted correspondences from images are refined to better match the surface meshes while maintaining the feature alignment learned in phase 2 while keeping the teacher network weights unchanged. The loss function for this phase is  $\mathcal{L} = \mathcal{L}_{EA} + \mathcal{L}_{PR}$  with  $\lambda_1 = 0$  and  $\lambda_2 = \lambda_3 = 1$ .

This comprehensive training strategy ensures optimal integration of surface representation based shape prior, for teaching the image encoder to learn representative shape features. During inference on testing samples, correspondences can be directly obtained from an image using the image encoder and the implicit field decoder. Additionally, to enhance the robustness of the surface autoencoder, we introduce vertex denoising as a data augmentation strategy during training. This is achieved by adding jitter to the input vertex positions, encouraging the autoencoder to learn to accurately denoise and reconstruct the original mesh vertices. The same data augmentation strategy can also be extended to point cloud data.

## 4 Datasets and Evaluation

### 4.1 Datasets

We select the left atrium and liver datasets for our experiments as they display highly variable shapes, which pose significant challenges for correspondence prediction tasks.

**Left Atrium Dataset (LA):** The dataset comprises 923 anonymized Late Gadolinium Enhancement (LGE) MRIs obtained from distinct patients and were manually segmented by cardiovascular medicine experts. The images were manually segmented at the University of Utah Division of Cardiovascular Medicine, the endocardium wall was used to cut off pulmonary veins. They have a spatial resolution of  $0.65 \times 0.65 \times 2.5mm^3$ , with the endocardial wall serving as the boundary for the pulmonary veins. Following segmentation, the images were cropped around the region of interest and downsampled by a factor of 0.8 to effectively manage memory usage, resulting in input images of size  $166 \times 120 \times 125$ .

**AbdomenCT-1K Liver Data:** The dataset [21] consists of CT scans and segmentations of four abdominal organs, including the liver, kidney, spleen, and pancreas. This dataset comprises 1132 3D CT scans sourced from various public datasets with segmentation verified and refined by experienced radiologists. We use this dataset’s CT scans and corresponding liver segmentations for the experiments. The CT scans have resolutions of  $512 \times 512$  pixels with varying pixel sizes and slice thicknesses between 1.25-5 mm. We visually assess the quality of the images and segmentations and utilize 833 samples. The images were cropped around the region of interest with the help of the segmentations and downsampled by a factor of 3.5 to manage memory usage effectively. The downsampled volume size is  $144 \times 156 \times 115$  with isotropic voxel spacing of 2 mm.

## 4.2 Models for Comparison

We compare the proposed model against the following:

1. **DeepSSM** [6] is a leading supervised model for predicting correspondence points from 3D image volumes. This method necessitates an optimized PDM for training, where each training instance consists of an image-correspondence pair. We utilize the correspondence supervised version of DeepSSM that uses a fixed decoder initialized with PCA basis and mean shape and trained on mean squared error (MSE) loss between predicted and ground truth correspondences.
2. **TL-DeepSSM** [6] is a variant of DeepSSM designed to overcome limitations associated with PCA usage. However, like DeepSSM, TL-DeepSSM is a supervised approach requiring optimized PDM and image pairs. The TL-variant network architecture [15] consists of a correspondence autoencoder and a T-flank network for image feature extraction. The network is trained using MSE between predicted and ground truth correspondences for the autoencoder and latent space MSE between the correspondence features and image features.
3. **Baseline** is introduced to demonstrate the effectiveness of introducing shape prior for the image based task. This model consists of an image encoder and an implicit field decoder trained end-to-end to predict correspondences by minimizing the Chamfer distance between the predicted correspondences from images and the mesh vertices.

### 4.3 Metrics

This section describes the metrics used to assess the performance of the quality of the shape models. Since SCorP does not use the ground truth PDM for training, we exclude the root mean square error metric, which is used in DeepSSM [7] and TL-DeepSSM [7,6].

1. **Chamfer distance (CD)** measures the average distance from each point in one set ( $C_j$ ) to its nearest neighbor in the other set ( $V_j$ ) and vice versa, providing a bidirectional measure of dissimilarity between two point sets.
2. **Point-to-mesh distance (P2M)** is the sum of point-to-mesh face distance and face-to-point distance for the predicted correspondences  $C_j$  and the mesh faces defined using vertices and edges ( $V_j, E_j$ ).
3. **Surface-to-surface (S2S) distance** is calculated between the original surface mesh and generated mesh from predicted correspondences. To obtain the reconstructed mesh, we match the correspondences to the mean shape and apply the warp between the points to its mesh.
4. **SSM Metrics:** Three statistical metrics are used to assess SSM correspondence [23]. **Compactness** refers to representing the training data distribution with minimal parameters, measured by the number of PCA modes needed to capture 95% of variation in correspondence points. **Generalization** evaluates how well the SSM extrapolates from training to unseen examples, gauged by the reconstruction error (L2) between held-out and training SSM-reconstructed correspondence points. **Specificity** measures the SSM’s ability to generate valid instances of the trained shape class, quantified by the average distance between sampled SSM correspondences and the nearest existing training correspondences.

### 4.4 Experimental Setup

For both datasets, we employ train/test/validation splits of 80%/10%/10%. We utilize ShapeWorks [9], an open-source shape modeling package, to process images and segmentations (align, crop, binarize segmentations, factor out scale and rotation) and generate surface meshes with 5000 vertices from the segmentations. Additionally, we use ShapeWorks [9] to generate the ground truth PDM and follow all prescribed procedures to obtain the required data for training DeepSSM and TL-DeepSSM [6,7]. We use the code and hyperparameters provided by the authors of DeepSSM and TL-DeepSSM [6] to train the models. The PDM is generated with 1024 correspondence particles, sufficient to capture the complex organ shapes.

We use the medoid shape of each dataset with 1024 correspondences as the template for the implicit field decoder. The medoid shape is identified using the surface-to-surface distances of meshes. This template remains consistent across all three training phases. We employ the Adam optimizer with a fixed learning rate of 0.00001 and continue training until convergence, determined through validation evaluation. Convergence is reached when the validation CD does not improve for 200 epochs. The models resulting from the epoch with the best validation CD are chosen for evaluation. Additionally, we set the weighting term

of the vertex reconstruction term  $\alpha$  to 0.001 (eq. 1) for all experiments. During the mesh branch training, a random jitter with a standard deviation of 1% of the maximum vertex size is added as data augmentation. The hyperparameters are identified via tuning for the validation set performance. The source code is available at [https://github.com/iyerkrithika21/SCorP\\_MIUA2024](https://github.com/iyerkrithika21/SCorP_MIUA2024).

To ensure a fair comparison, we maintain consistency by employing the same architecture for the image encoder and replicating the same image normalization steps used in DeepSSM and TL-DeepSSM [7,6] across all experiments. This approach mitigates potential biases arising from architectural variations, facilitating an accurate assessment of performance differences.

## 5 Results

Fig. 3 shows an overview of the metrics for the held-out test samples of the liver and LA datasets. The baseline method, trained without the mesh-informed shape prior while using the same inference architecture, demonstrates inferior performance for all metrics across the two datasets. This finding emphasizes the critical role of leveraging shape prior information to enhance image-based SSM prediction tasks, especially in the absence of supervision. The proposed model, SCorP, outperforms the other methods in terms of CD. SCorP also performs better with respect to P2M and S2S distances for the LA dataset. On the liver dataset, SCorP exhibits competitive performance with P2M and S2S distances. SCorP provides the best compactness for both datasets suggesting strong correspondence and showcases comparable specificity and generalization.

Furthermore, Fig. 4 and Fig. 5 depict the top four Principal Component Analysis (PCA) modes of variation identified by SCorP, DeepSSM, and TL-DeepSSM for the LA and liver datasets, respectively. Despite being an unsupervised method, SCorP demonstrates competitive performance in identifying modes of variation. Additionally, SCorP shows detailed and smoother variations as compared to the other methods, which are highlighted with boxes in Fig. 4 and Fig. 5.

We examined the worst and median-performing samples in terms of the P2M distance for all three methods and discovered a substantial overlap among them, suggesting similar success and failure modes. We overlaid the true surface mesh for two median cases and two worst-performing samples with the correspondence-level P2M distances, as depicted in Fig. 6. Additionally, we analyze the corresponding image slices to gain insights into the performance discrepancies. For the LA dataset, comparing Fig. 6.A and Fig. 6.B reveals the significant impact of image quality on the performance of all three methods. Additionally, a notable observation is the deviation of the shape of the worst-case sample in Fig. 6.B from the population mean (see mean shape in Fig. 4). Similarly, for the liver dataset, comparing Fig. 6.C and Fig. 6.D highlights clear distinctions in image quality between median and worst-case scenarios. The image slices corresponding to the worst P2M distance exhibit poor contrast and an unclear picture of the liver shape, posing challenges for the image encoder. Notably, when exam-

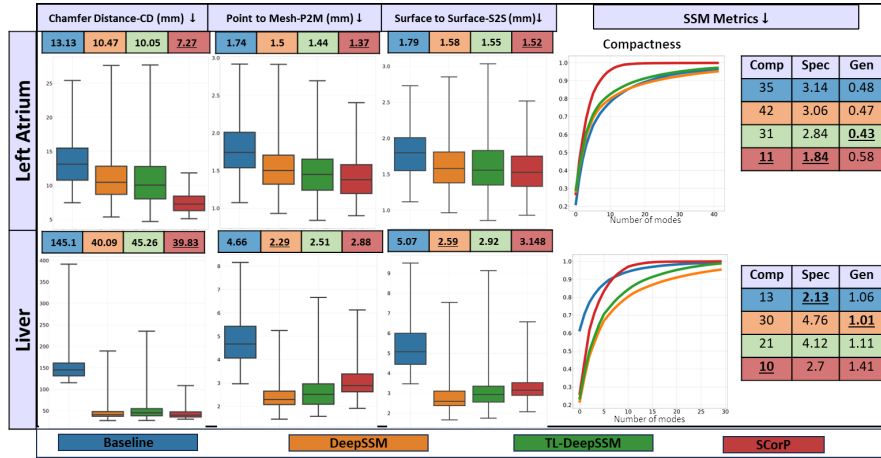


Fig. 3: **Performance metrics** Boxplots illustrating the distribution of performance metrics, with mean values displayed above each plot, for the held-out test samples from the LA and liver datasets. Compactness plots illustrate the cumulative population variation captured by PCA modes, where a larger area under the curve indicates a more compact model. The best metrics are **highlighted** in the figure. Comp = Compactness, Spec = Specificity, Gen = Generalization.

in Fig. 6, we observe that SCorP performs comparably with DeepSSM and TL-DeepSSM for median cases. However, for worst-case P2M scenarios in the first row of Fig. 6.C and Fig. 6.D, SCorP demonstrates superior performance, producing better correspondences for the same samples compared to DeepSSM and TL-DeepSSM.

### 5.1 Ablation Experiments

*Impact of training sample size:* We also analyze the robustness of all methods across different training dataset sizes (15%, 20%, 40%, 80%, 100%), which yields valuable insights. Fig. 7 illustrates clear trends in mean performance metrics and their standard deviations across various methods at each dataset size. As expected, expanding the training dataset size leads to improved performance across all metrics for all models. Interestingly, even with smaller dataset sizes, SCorP consistently outperforms DeepSSM and TL-DeepSSM, indicating its robustness and superior generalization ability. One contributing factor to this trend is that the SCorP does not rely on an optimized PDM during training, unlike DeepSSM and TL-DeepSSM, which depend on the optimized PDM, imposing stronger linearity constraints thereby limiting generalization, particularly with smaller training datasets. SCorP exhibits greater flexibility and adaptability, allowing it to achieve competitive performance even with limited training data.

*Point clouds for surface representation:* To demonstrate the versatility of SCorP across different surface representation formats, we experimented using the LA dataset, employing point clouds sampled from the meshes to encode the feature

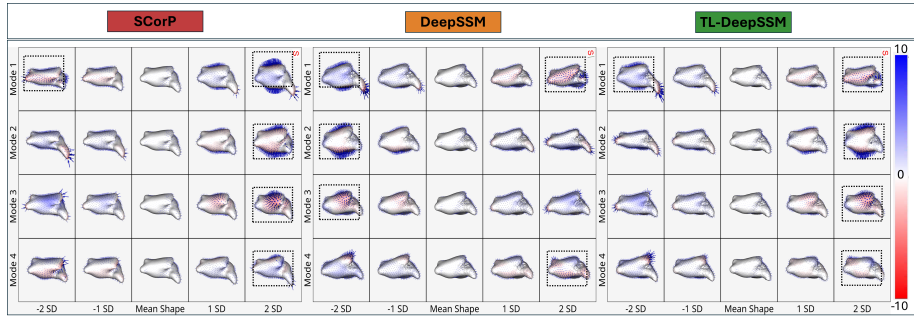


Fig. 4: **PCA modes of variations:** The first four modes of variations of the LA dataset identified by SCorP, DeepSSM, and TL-DeepSSM [7,6]. The color map and arrows show the signed distance and direction from the mean shape. SCorP shows detailed and smoother variations as compared to the other methods, which are highlighted with boxes.

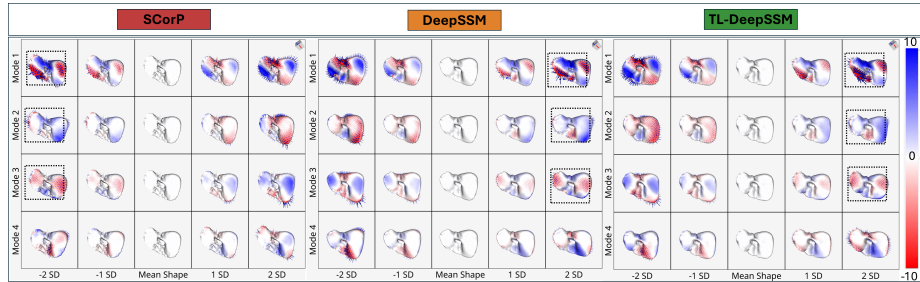


Fig. 5: **PCA modes of variations:** The first four modes of variations of the liver dataset identified by SCorP, DeepSSM, and TL-DeepSSM [7,6]. The color map and arrows show the signed distance and direction from the mean shape. SCorP shows detailed and smoother variations as compared to the other methods, which are highlighted with boxes.

vector. In this setup, we utilized the Euclidean distance for k-nearest neighbor calculation in the initial layer of the DGCNN mesh encoder [29]. Following the training steps outlined in Section 3.3, the model exhibited performance similar to the best-performing model from Fig. 3 with the following statistics: **CD**  $7.512 \pm 1.72$ , **P2M**  $1.435 \pm 0.326$ , and **S2S**  $1.56 \pm 0.348$ . This highlights that our model is agnostic to the underlying surface representation, enhancing its generalization and usability compared to its counterparts.

## 6 Limitations and Future Work

Given the pivotal role of SSM in diagnostic clinical support systems, it is critical to address the limitations of SCorP. The model currently requires the cohort of images and shapes to be aligned. Relaxing this requirement through developing robust alignment algorithms or exploring alignment-free methods can broaden the usability of SCorP across various datasets and clinical scenarios.



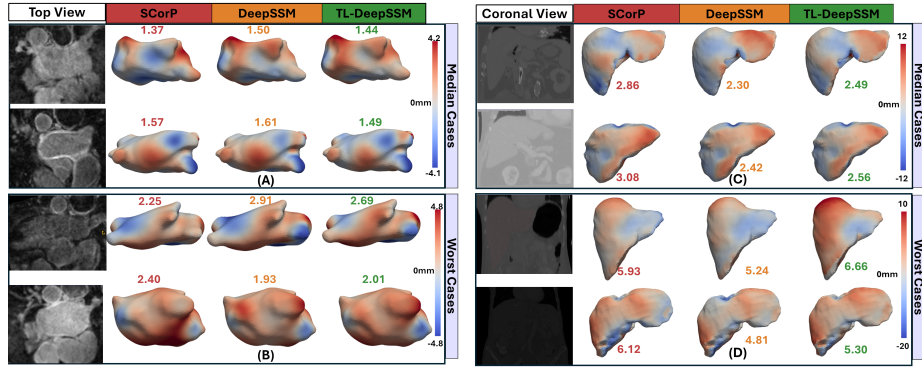


Fig. 6: The volumetric image slices of representative samples (worst and median cases) for LA (A and B) and liver (C and D) datasets and all models. The ground truth meshes of the representative samples with a distance map overlay with the correspondence-wise P2M distances for the respective models. The numbers above each sample represent the absolute average P2M distance of the sample. All the models have similar modes of failure and success, and the performances are affected by image quality and the degree of shape outlier.

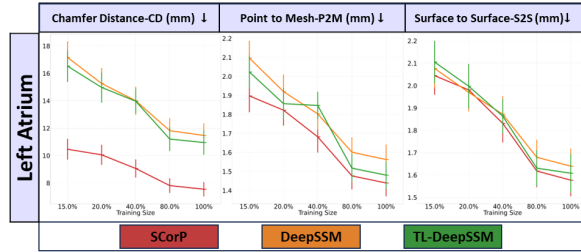


Fig. 7: **Impact of training dataset sizes** The plot illustrates the mean and standard deviation of the performance metrics for all methods at 15%, 20%, 40%, 80%, 100%) training dataset size. SCorP consistency outperforms DeepSSM and TL-DeepSSM and proves to be robust to the dataset size.

Furthermore, expanding the framework’s capabilities to accommodate diverse data types (sparse slices, orthogonal view slices, and radiography data) and incorporating data augmentation schemes (similar to ADASSM [19]) holds immense promise for broadening its applicability. Additionally, integrating uncertainty quantification methods to evaluate prediction confidence would enhance result interpretability, further advancing the utility of SCorP in clinical settings.

## 7 Conclusion

The proposed framework, SCorP presents a novel approach to inferring correspondences directly from raw images without needing a pre-optimized shape

model. By integrating prior shape information from surface representations (meshes, point clouds, binary volumes), SCorP achieves superior performance compared to traditional and state-of-the-art deep learning methods with less supervision. The three-phase training strategy ensures effective integration of shape statistics-informed priors, guiding the image encoder to learn representative shape features for the correspondence prediction task. Furthermore, SCorP demonstrates robustness across varying training dataset sizes, highlighting its versatility and applicability in different scenarios. Overall, SCorP improves upon existing methods by streamlining the PDM generation process, which increases the feasibility of using shape models for research and applications in medical imaging, computer-aided diagnosis, and beyond.

## 8 Acknowledgements

This work was supported by the National Institutes of Health under grant numbers NIBIB-U24EB029011, NIAMS-R01AR076120, and NHLBI-R01HL135568. 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.: Fully bayesian vib-deepssm. In: International Conference on Medical Image Computing and Computer-Assisted Intervention. pp. 346–356. Springer (2023)
5. Aziz, A.Z.B., Adams, J., Elhabian, S.: Progressive deepssm: Training methodology for image-to-shape deep models. In: International Workshop on Shape in Medical Imaging. pp. 157–172. Springer (2023)
6. 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)
7. 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: International Workshop, ShapeMI 2018, Held in Conjunction with MICCAI 2018, Granada, Spain, September 20, 2018, Proceedings. pp. 244–257. Springer (2018)
8. Borotikar, B., Mutsvangwa, T.E., Elhabian, S.Y., Audenaert, E.A.: Statistical model-based computational biomechanics: applications in joints and internal organs. *Frontiers in Bioengineering and Biotechnology* **11**, 1232464 (2023)

9. 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)
10. Cerrolaza, J.J., Picazo, M.L., Humbert, L., Sato, Y., Rueckert, D., Ballester, M.Á.G., Linguraru, M.G.: Computational anatomy for multi-organ analysis in medical imaging: A review. *Medical Image Analysis* **56**, 44–67 (2019)
11. Chen, Z.: Im-net: Learning implicit fields for generative shape modeling (2019)
12. Davies, R.H.: *Learning shape: optimal models for analysing natural variability*. The University of Manchester (United Kingdom) (2002)
13. Durrleman, S., Prastawa, M., Charon, N., Korenberg, J.R., Joshi, S., Gerig, G., Trounev, A.: Morphometry of anatomical shape complexes with dense deformations and sparse parameters. *NeuroImage* **101**, 35–49 (2014)
14. Friedrich, P., Wolleb, J., Bieder, F., Thieringer, F.M., Cattin, P.C.: Point cloud diffusion models for automatic implant generation. In: *International Conference on Medical Image Computing and Computer-Assisted Intervention*. pp. 112–122. Springer (2023)
15. Girdhar, R., Fouhey, D.F., Rodriguez, M., Gupta, A.: Learning a predictable and generative vector representation for objects. In: *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part VI 14*. pp. 484–499. Springer (2016)
16. Heitz, G., Rohlfing, T., Maurer Jr, C.R.: Statistical shape model generation using nonrigid deformation of a template mesh. In: *Medical Imaging 2005: Image Processing*. vol. 5747, pp. 1411–1421. SPIE (2005)
17. 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)
18. Jiang, C., Huang, J., Tagliasacchi, A., Guibas, L.J.: Shapeflow: Learnable deformation flows among 3d shapes. *Advances in Neural Information Processing Systems* **33**, 9745–9757 (2020)
19. Karanam, M.S.T., Kataria, T., Iyer, K., Elhabian, S.Y.: Adassm: Adversarial data augmentation in statistical shape models from images. In: *International Workshop on Shape in Medical Imaging*. pp. 90–104. Springer (2023)
20. Lüdke, D., Amiranashvili, T., Ambellan, F., Ezhov, I., Menze, B.H., Zachow, S.: Landmark-free statistical shape modeling via neural flow deformations. In: *Medical Image Computing and Computer Assisted Intervention–MICCAI 2022: 25th International Conference, Singapore, September 18–22, 2022, Proceedings, Part II*. pp. 453–463. Springer (2022)
21. 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.: Abdomenct-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>
22. Mori, N., Mugikura, S., Endo, T., Endo, H., Oguma, Y., Li, L., Ito, A., Watanabe, M., Kanamori, M., Tominaga, T., et al.: Principal component analysis of texture features for grading of meningioma: not effective from the peritumoral area but effective from the tumor area. *Neuroradiology* **65**(2), 257–274 (2023)
23. 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)

24. Riordan, E., Yung, A., Cheng, K., Lim, L., Clark, J., Rtshiladze, M., Ch'ng, S.: Modeling methods in craniofacial virtual surgical planning. *Journal of Craniofacial Surgery* **34**(4), 1191–1198 (2023)
25. Samson, C., Blanc-Féraud, L., Aubert, G., Zerubia, J.: A level set model for image classification. *International journal of computer vision* **40**(3), 187–197 (2000)
26. 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* (1071), 242 (2006)
27. Tufegdžic, M., Trajanovic, M.D.: Building 3d surface model of the human hip bone from 2d radiographic images using parameter-based approach. In: *Personalized Orthopedics: Contributions and Applications of Biomedical Engineering*, pp. 147–181. Springer (2022)
28. Ukey, J., Elhabian, S.: Localization-aware deep learning framework for statistical shape modeling directly from images. In: *Medical Imaging with Deep Learning* (2023)
29. Wang, Y., Sun, Y., Liu, Z., Sarma, S.E., Bronstein, M.M., Solomon, J.M.: Dynamic graph cnn for learning on point clouds. *Acm Transactions On Graphics (tog)* **38**(5), 1–12 (2019)
30. Xu, H., Elhabian, S.Y.: Image2ssm: Reimagining statistical shape models from images with radial basis functions. In: *International Conference on Medical Image Computing and Computer-Assisted Intervention*. pp. 508–517. Springer (2023)
31. Zhu, C., Dev, H., Sharbatdaran, A., He, X., Shimonov, D., Chevalier, J.M., Blumenfeld, J.D., Wang, Y., Teichman, K., Shih, G., et al.: Clinical quality control of mri total kidney volume measurements in autosomal dominant polycystic kidney disease. *Tomography* **9**(4), 1341–1355 (2023)

## A Appendix

### A.1 Hyperparamters

All models were trained on NVIDIA GeForce RTX 2080 Ti.

Parameter	Description	Value
B	Batch size	6
LR	Learning rate	$1e^{-5}$
M	Number of correspondences	1024
ES	Early stopping patience epochs	200

Table 1: Hyperparameters shared by all models

### A.2 Architecture

1. **Image encoder:** The encoder architecture utilizes Conv2d layers with  $5 \times 5$  filters and the following numbers of filters: [12, 24, 48, 96, 192]. After each Conv2d layer, batch normalization and ReLU activation functions are applied. Max pooling layers are incorporated to reduce spatial dimensions. The

Parameter	Description	Value
L	Latent dimension for SCorP	256
K	Size of neighbourhood for EdgeConv	20 (LA), 27 (Liver)
NV	Number of vertices in the mesh	5000

Table 2: Hyperparameters for SCorP

feature maps are then flattened and passed to the fully connected layers. The fully connected (FC) layer stack consists of linear layers with different input and output feature dimensions: [193536- > 384], [384- > 96], [96- > 256]. Each linear layer is followed by a Parametric ReLU (PReLU) activation function.

2. **Surface Autoencoder:** We use the *DGCNN\_semseg\_s3dis* model from the original DGCNN Github repository.
3. **IM-Net:** We use the original implementation of IM-Net from the Github repository.

### A.3 SSM Metrics

1. Compactness: We quantify compactness as the number of PCA modes that are required to capture 95% of the total variation in the output training cohort correspondence points.
2. Specificity: We quantify specificity by randomly generating  $J$  samples from the shape space using the eigenvectors and eigenvalues that capture 95% variability of the training cohort. Specificity is computed as the average squared Euclidean distance between these generated samples and their closest training sample.

$$S = \sum_{\mathcal{C} \in \mathcal{C}_{generated}} \|\mathcal{C} - \mathcal{C}_{train}\|^2$$

3. Generalization: We quantify generalization by assessing the average approximation errors across a set of unseen instances. Generalization is defined as the mean approximation errors between the original unseen shape instance and reconstruction of the shape constructed using the training cohort PCA eigenvalues and vectors that preserve 95% variability.

$$G = \sum_{j=1}^U \|\mathcal{C}_j - \hat{\mathcal{C}}_j\|_2^2 \text{ for } J \text{ unseen shapes.}$$