

3D-LMNet: Latent Embedding Matching for Accurate and Diverse 3D Point Cloud Reconstruction from a Single Image

Priyanka Mandikal*
priyanka.mandikal@gmail.com

K L Navaneet*
navaneetl@iisc.ac.in

Mayank Agarwal*
mayankgrwl97@gmail.com

H. Venkatesh Babu
venky@iisc.ac.in

Video Analytics Lab,
Department of Computational and
Data Sciences,
Indian Institute of Science,
Bangalore, India

Abstract

3D reconstruction from single view images is an ill-posed problem. Inferring the hidden regions from self-occluded images is both challenging and ambiguous. We propose a two-pronged approach to address these issues. To better incorporate the data prior and generate meaningful reconstructions, we propose 3D-LMNet, a latent embedding matching approach for 3D reconstruction. We first train a 3D point cloud auto-encoder and then learn a mapping from the 2D image to the corresponding learnt embedding. To tackle the issue of uncertainty in the reconstruction, we predict multiple reconstructions that are consistent with the input view. This is achieved by learning a probabilistic latent space with a novel view-specific ‘diversity loss’. Thorough quantitative and qualitative analysis is performed to highlight the significance of the proposed approach. We outperform state-of-the-art approaches on the task of single-view 3D reconstruction on both real and synthetic datasets while generating multiple plausible reconstructions, demonstrating the generalizability and utility of our approach.

1 Introduction



Figure 1: Single-view reconstructions for unambiguous and ambiguous input views.

* equal contribution

Humans can infer the structure of a scene and the shapes of objects within it from limited information. Even for regions that are highly occluded, we are able to guess a number of plausible shapes that could complete the object. Our ability to directly perceive the 3D structure from limited 2D information arises from a strong prior about shapes and geometries that we are familiar with. This ability is central to our perception of the world and the manipulation of objects within it.

Extending the above idea to machines, the ability to infer the 3D structures from single-view images has far-reaching applications in the field of robotics and perception, in tasks such as robot grasping, object manipulation, etc. However, the task is particularly challenging due to the inherent ambiguity that exists in the reconstructions of occluded images. While the existing data-driven approaches capture the semantic information present in the image to accurately reconstruct corresponding 3D models, it is unreasonable to expect them to predict a single deterministic output for an ambiguous input. An ideal machine would produce multiple solutions when there is uncertainty in the input, while obtaining a deterministic output for images with adequate information (Fig. 1).

With the recent advances of deep learning, the problem of 3D reconstruction has largely been tackled with the help of 3D-CNNs that generate a voxelized 3D occupancy grid. However, this representation suffers from sparsity of information, since most of the information needed to perceive the 3D structure is provided by the surface voxels. 3D CNNs are also compute heavy and add considerable overhead during training and inference. To overcome the drawbacks of the voxel representation, recent works have focused on designing neural network architectures and loss formulations to process and predict 3D point clouds [8, 15, 16], which consist of points being sampled uniformly on the object’s surface. The information-rich encoding and compute-friendly architectures makes it an ideal candidate for 3D shape generation and reconstruction tasks. Hence, we consider point clouds as our 3D representations.

In this work, we seek to answer two important questions in the task of single-view reconstruction **(1)** Given a two-dimensional image of an object, what is an effective way of inferring an accurate 3D point cloud representation of it? **(2)** When the input image is highly occluded, how do we equip the network to generate a set of plausible 3D shapes that are consistent with the input image? We achieve the former by first learning a strong prior over all possible 3D shapes with the help of a 3D point cloud auto-encoder. We then train an image encoder to map the input image to this learnt latent space. To address the latter issue, we propose a mechanism to learn a probabilistic distribution in the latent space that is capable of generating multiple plausible outputs from possibly ambiguous input views.

In summary, our contributions in this work are as follows:

- We propose a latent-embedding matching setup called 3D-LMNet, to demonstrate the importance of learning a good prior over 3D point clouds for effectively transferring knowledge from the 3D to 2D domain for the task of single-view reconstruction. We thoroughly evaluate various ways of mapping to a learnt 3D latent space.
- We present a technique to generate multiple plausible 3D shapes from a single input image to tackle the issue of ambiguous ground truths, and empirically evaluate the effectiveness of this strategy in generating diverse predictions for ambiguous views.
- We evaluate 3D-LMNet on real data and demonstrate the generalizability of our approach, which significantly outperforms the state-of-art reconstruction methods for the task of single-view reconstruction.

2 Related Work

3D Reconstruction

With the advent of deep neural network architectures in 2D image generation tasks, the power of convolutional neural nets have been directly transferred to the 3D domain using 3D CNNs. A number of works have revolved around generating voxelized output representations [1, 9, 23, 24]. Giridhar *et al.* [9] learnt a joint embedding of 3D voxel shapes and their corresponding 2D images. While the focus of [9] was to learn a vector representation that is generative and predictable at the same time, our aim is to address the problem of transferring the knowledge learnt in 3D to the 2D domain specifically for the task of single-view reconstruction. Additionally, we tackle the rather under-addressed problem of generating multiple plausible outputs that satisfy the given input image. Wu *et al.* [23] used adversarial training in a variational setup for learning more realistic generations. Choy *et al.* [9] trained a recurrent neural network to encode information from more than one input views. Works such as [22, 25] explore ways to reconstruct 3D shapes from 2D projections such as silhouettes and depth maps. Apart from reconstructing shapes from scratch, other reconstruction tasks such as shape completion [8, 19] and shape deformation [27] have also been studied in the voxel domain. But the compute overhead and sparsity of information in voxel formats inspired lines of work that abstracted volumetric information into smaller number of units with the help of the octree data structure [11, 18, 21].

More recently, Fan *et al.* [8], introduced frameworks and loss formulations tailored for generating unordered point clouds, and achieved single-view 3D reconstruction results outperforming the volumetric state-of-art approaches [9]. While [8] directly predicts the 3D point cloud from 2D images, our approach stresses the importance of first learning a good 3D latent space of point clouds before mapping the 2D images to it. Lin *et al.* [12] generated point clouds by fusing depth images and refined them using a projection module. Apart from single-view reconstruction, there is active research in other areas of point cloud analysis including processing [6, 17], upsampling [26], deformation [13], and generation [10].

Generating multiple plausible outputs

While multiple correct reconstructions can exist for a single input image, most prior works predict deterministic outputs regardless of the information that is available. Rezende *et al.* [17] and Fan *et al.* [8] tackle the problem by training a conditional variational auto-encoder [9, 11] on 3D shapes conditioned on the input image. In [8], an alternative approach of inducing randomness into the model at the input stage is considered. In 3D-LMNet, we introduce a training regime comprising of sampling a probabilistic latent variable, and optimizing a novel view-specific loss function. Our reconstructions exhibit greater semantic diversity and effectively model the view-specific uncertainty present in the data distribution.

3 Approach

Our training pipeline consists of two stages as outlined in Fig. 2. In the first stage, we learn a latent space $\mathcal{Z} \subseteq \mathbb{R}^k$ of 3D point clouds by training a point cloud auto-encoder (E_p, D_p). In the second stage, we train an image encoder E_i to map the 2D image to this learnt latent space \mathcal{Z} . A variant of stage two consists of mapping to \mathcal{Z} in a probabilistic manner so as to infer multiple possible predictions for a single input image during test time. Each of the components is described below in detail.

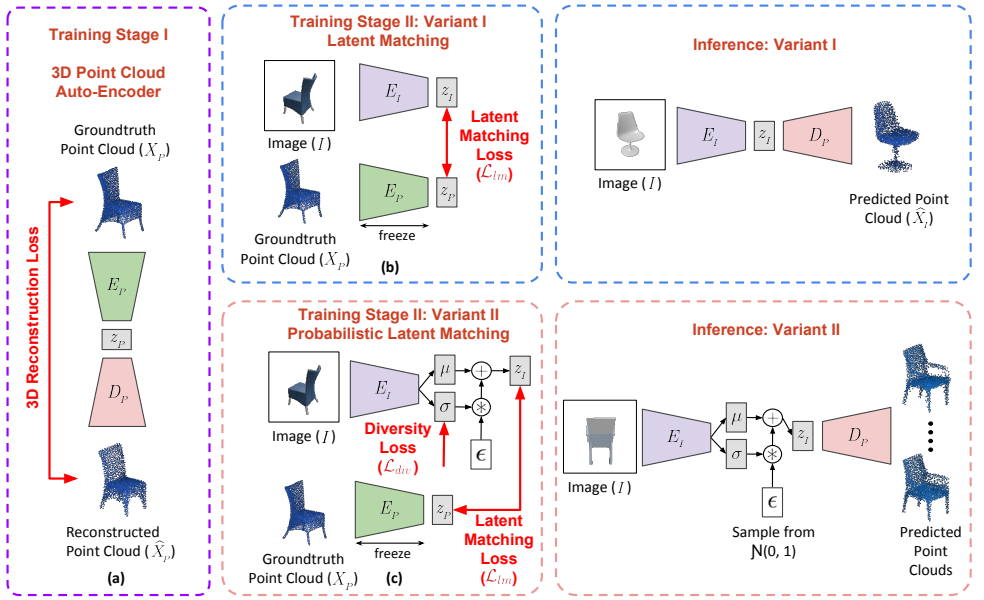


Figure 2: Overview of 3D-LMNet. The training pipeline consists of two stages. In Stage I, we learn a latent space \mathcal{Z} for 3D point clouds by training a point cloud auto-encoder (E_p , D_p). In Stage II, we train an image encoder E_i to map the 2D images to this learnt \mathcal{Z} . In a variant of Stage II, we map to \mathcal{Z} in a probabilistic manner so as to infer multiple plausible predictions for a single input image during inference.

3.1 3D Point Cloud Auto-Encoder

Our goal is to learn a strong prior over the 3D point clouds in the dataset. For this purpose, we train an encoder-decoder network (E_p , D_p) that takes in a ground truth point cloud $X_p \in \mathbb{R}^{N \times 3}$ and outputs a reconstructed point cloud $\hat{X}_p \in \mathbb{R}^{N \times 3}$, where N is the number of points in the point cloud (Fig. 2a). Since a point cloud is an unordered form of representation, we need a network architecture that is invariant to the relative ordering of input points. To enforce this, we choose the architecture of E_p based on PointNet [15], consisting of 1D convolutional layers acting independently on every point in the point cloud X_p . To achieve order-invariance of point features in the latent space, we apply the maxpool symmetry function to obtain a bottleneck \mathcal{Z} of dimension k . The decoder consists of fully-connected layers operating on \mathcal{Z} to produce the reconstructed point cloud \hat{X}_p . Since the loss function for optimization also needs to be order-invariant, Chamfer distance between X_p and \hat{X}_p is chosen as the reconstruction loss. The loss function is defined as:

$$L_{rec} = d_{Chamfer}(X_p, \hat{X}_p) = \sum_{x \in X_p} \min_{\hat{x} \in \hat{X}_p} \|x - \hat{x}\|_2^2 + \sum_{\hat{x} \in \hat{X}_p} \min_{x \in X_p} \|x - \hat{x}\|_2^2 \quad (1)$$

Once the auto-encoder is trained, the next stage consists of training an image encoder to map to this learnt embedding space.

3.2 Latent Matching

In this stage, we aim to effectively transfer the knowledge learnt in the 3D point cloud domain to the 2D image domain. We train an image encoder E_I that takes in an input image I and outputs a latent image vector z_I of dimension k (Variant I, Fig. 2b). There are two ways of achieving the 3D-to-2D knowledge transfer:

- (1) **Matching the reconstructions:** We pass \hat{z}_I through the pre-trained point cloud decoder D_p to get the predicted point cloud \hat{X}_I . The parameters of D_p are not updated during this step. Chamfer distance between \hat{X}_I and X_p is used as the loss function for optimization. We refer to this variant as "3D-LMNet-Chamfer" in the evaluation section (Sec. 4).
- (2) **Matching vectors in the latent \mathcal{Z} space:** The latent representations of the image and corresponding ground truth point cloud are matched. The error is computed between the predicted z_I and the ground truth z_p , obtained via passing X_p through the pre-trained point cloud encoder E_p (Fig. 2b). The parameters of E_p are not updated during this step. For the latent loss \mathcal{L}_{lm} , we experiment with the squared euclidean error ($\mathcal{L}_2(z_I - z_p) = \|z_I - z_p\|_2^2$) and the least absolute error ($\mathcal{L}_1(z_I - z_p) = |z_I - z_p|$) for matching the latent vectors. We refer to these two variants as "3D-LMNet- \mathcal{L}_2 " and "3DLMNet- \mathcal{L}_1 " in the evaluation section (Sec. 4). During inference, we obtain the predicted point cloud by passing the image through E_I followed by D_p .

In our experiments (detailed in Sec. 4.1), we find that alternative two i.e. matching latent vectors provides substantial improvement over optimizing for the reconstruction loss.

3.3 Generating Multiple Plausible Outputs

We propose to handle the uncertainty in predictions by learning a probabilistic distribution in the latent space \mathcal{Z} . For every input image I in the dataset, there are multiple settings of the latent variables z for which the model should predict an output that is consistent with I . To allow the network to make probabilistic predictions, we formulate the latent representation z_I of a specific input image I_1 to be a Gaussian random variable, i.e. $z_I \sim \mathcal{N}(\mu, \sigma^2)$ (Variant II, Fig. 2c). Similar to Variational Auto-Encoders (VAE) [14], we use the "reparameterization trick" to handle stochasticity in the network. The image encoder predicts the mean μ and standard deviation σ of the distribution, and $\varepsilon \sim \mathcal{N}(0, 1)$ is sampled to obtain the latent vector as $z_I = \mu + \varepsilon\sigma$ (Fig. 2). However, unlike in the case of conventional VAEs, the mean of the distribution is unconstrained, while the variance is constrained such that meaningful and diverse reconstructions are obtained for a given input image.

A critical challenge is to obtain a model that can generate diverse but semantically meaningful predictions for an occluded view while retaining the visible semantics. Another challenge is to make highly confident predictions when the input view is informative. To accommodate this, we formulate a fast-decaying loss function that penalizes σ for being too far off from zero for unambiguous views, while giving it the liberty to explore the latent space for ambiguous views. We term this as the *diversity loss* and define it as follows:

$$\mathcal{L}_{div} = \left(\sigma - \eta e^{-\frac{(\phi_i - \phi_o)^2}{\delta^2}} \right)^2 \quad (2)$$

where, ϕ_i is the azimuth angle of the input image I , ϕ_o is the azimuth angle of maximum occlusion view, and δ determines the rate of decay. η controls the magnitude of standard

Metric	AE	Baseline	3D-LMNet Chamfer	3D-LMNet \mathcal{L}_2	3D-LMNet \mathcal{L}_1
Chamfer	4.46	5.78	5.99	5.54	5.40
EMD	6.53	9.20	7.82	7.20	7.00

Table 1: A comparison of the baseline and different variants of 3D-LMNet for the task of 3D reconstruction on ShapeNet [9]. All metrics are scaled by 100.

Metric	3D-LMNet Chamfer	3D-LMNet \mathcal{L}_2	3D-LMNet \mathcal{L}_1
\mathcal{L}_2	56.7	1.32	1.38
\mathcal{L}_1	14.02	1.34	1.29

Table 2: A comparison of latent matching errors for different variants of 3D-LMNet on ShapeNet [9]. All metrics are scaled by 0.01.

deviation σ . The above formulation can easily be extended to cases where multiple highly occluded views are present by considering a mixture of Gaussians.

The joint optimization loss function is a combination of the latent matching loss \mathcal{L}_{lm} and the diversity loss \mathcal{L}_{div} :

$$\mathcal{L} = \mathcal{L}_{lm} + \lambda \mathcal{L}_{div} \quad (3)$$

where λ is the weighing factor. During inference, the model is capable of generating diverse predictions when ε is varied. Note that pose information is not used during inference.

3.4 Implementation Details

In the point cloud auto-encoder, the encoder consists of five 1D convolutional layers with [64, 128, 128, 256] filters, ending with a bottleneck layer of dimension 512. We choose max-pool function as the symmetry operation. The decoder consists of three fully-connected layers of size [256, 256, $N \times 3$], where N is the number of points predicted by our network. We set N to be 2048 in all our experiments. We use the ReLU non-linearity and batch-normalization at all layers of the auto-encoder. The image encoder is a 2D convolutional neural network that maps the input image to the 512-dimensional latent vector. We use the Adam optimizer with a learning rate of 0.00005 and a minibatch size of 32. Network architectures for all components in our proposed framework are provided in the supplementary material. Codes are available at <https://github.com/val-iisc/3d-lmnet>.

4 Experiments

Dataset: We train all our networks on synthetic models from the ShapeNet [9] dataset. We use the same 80% – 20% train/test split provided by [9] consisting of models from 13 different categories, so as to be comparable with the previous works.

Evaluation Methodology: We report both the Chamfer Distance (Eqn. 1) as well as the Earth Mover’s Distance (or EMD) computed on 1024 randomly sampled points in all our

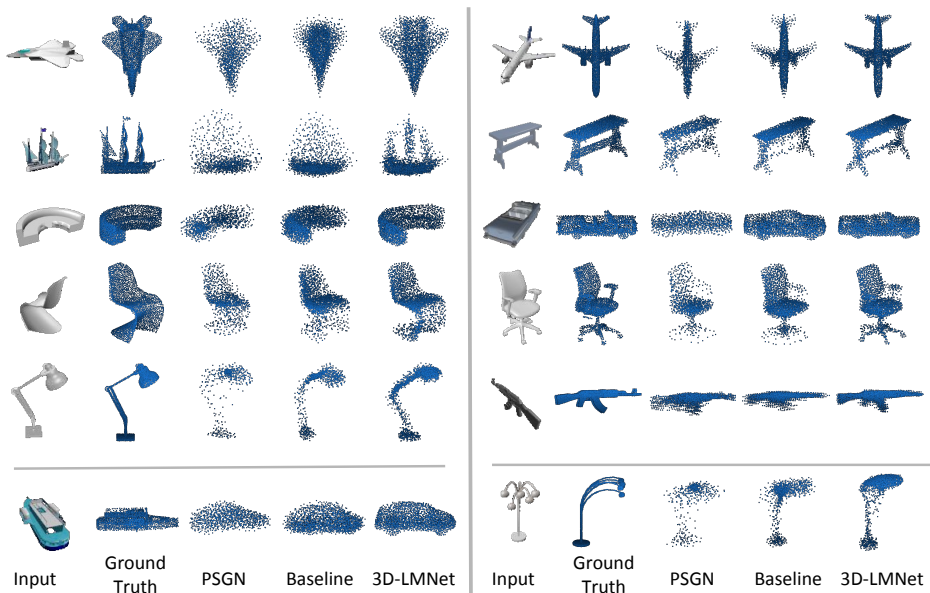


Figure 3: Qualitative results on ShapeNet [8]. Compared to PSGN [8] and the baseline, we are better able to capture the overall shape and finer details present in the input image. While clusters and outlier points are present in PSGN and baseline reconstructions, we obtain more uniformly distributed points. The bottom row presents failure cases for our approach. Note that PSGN predicts 1024 points, while the baseline and 3D-LMNet predict 2048 points.

evaluations. EMD between two point sets X_p and \hat{X}_p is given by:

$$d_{EMD}(X_p, \hat{X}_p) = \min_{\phi: X_p \rightarrow \hat{X}_p} \sum_{x \in X_p} \|x - \phi(x)\|_2 \quad (4)$$

where $\phi : X_p \rightarrow \widehat{X}_p$ is a bijection. We use an approximate. For computing the metrics, we renormalize both the ground truth and predicted point clouds within a bounding box of length 1 unit. Since PSGN [8] outputs are non-canonical, we align their predictions to the canonical ground truth frame by using pose metadata available in the evaluation datasets. Additionally, we apply the iterative closest point algorithm (ICP) [2] on the ground truth and predicted point clouds for finer alignment.

4.1 Empirical Evaluation on ShapeNet

We study the framework presented and evaluate each of the components in the training pipeline. To show the advantage of our latent matching procedure over direct 2D-3D training, we train a baseline which consists of an encoder-decoder network that is trained end-to-end, using reconstruction loss on the generated point cloud. We employ the same network architecture as the one used for latent matching experiments. To measure the performance of the baseline and all variants of our model, we use the validation split provided by [9] for reporting the Chamfer and EMD metrics. Table 1 shows the comparison between the baseline

and three variants of loss formulation in our latent matching setup. We also report the auto-encoder reconstruction scores, which serve as an upper bound on the performance of latent matching. We observe that the latent matching variants of 3D-LMNet outperform the baseline in both Chamfer and EMD. Amongst the variants, we observe that trivially optimizing for Chamfer loss leads to worse results, whereas training with losses directly operating on the latent space results in lower reconstruction errors. We also see that the \mathcal{L}_1 loss formulation performs better both in terms of Chamfer and EMD metrics. Additionally, we also report the latent matching errors for different variants of 3D-LMNet in Table 2. We observe that more accurate latent matching (characterized by lower \mathcal{L}_1 and \mathcal{L}_2 errors), results in lower reconstruction errors as well (Table 1). Category-wise metrics for all the variants are provided in the supplementary material.

4.2 Comparison with other methods on ShapeNet and Pix3D

We compare our 3D-LMNet- \mathcal{L}_1 model with PSGN [8] on the synthetic ShapeNet dataset [9] and the more recent Pix3D dataset [20] to test for generalizability on real world data. Since [8] establishes that point cloud based approach significantly outperforms the state-of-art voxel based approaches, we do not show any comparison against them.

ShapeNet Table 4 shows the comparison between PSGN [8], the baseline and our \mathcal{L}_1 latent matching variant on the validation split provided by [9]. We outperform PSGN in 8 out of 13 categories in Chamfer and 10 out of 13 categories in the EMD metric, while also having lower overall mean scores. It is worth noting that we achieve state-of-the-art performance in both metrics despite the fact that our network has half the number of trainable parameters in comparison to PSGN, while predicting point clouds with double the resolution. A lower EMD score also correlates with better visual quality and encourages points to lie closer to the surface [0, 26]. Qualitative comparison is shown in Fig. 3. Compared to PSGN [8] and the baseline, we are better able to capture the overall shape and finer details present in the input image. Note that both the other methods have clustered points and outlier points while our reconstructions are more uniformly distributed. We also present two failure cases of our approach in Fig. 3 (bottom row). Interestingly, we observe that in some cases, latent matching incorrectly maps an image to a similar looking object of different category, leading to good-looking but incorrect reconstructions. Fig. 3 shows a vessel being mapped to a car of similar shape. Another common failure case is the absence of finer details in the reconstructions. However, other approaches also have this drawback.

Pix3D For testing the generalizability of our approach on real-world datasets, we evaluate the performance of our method on the Pix3D dataset [20]. It consists of a large collection of real images and their corresponding metadata such as masks, ground truth CAD models and pose. We evaluate our trained model on categories that co-occur in the synthetic training set and exclude images having occlusion and truncation from the test set, as is done in the original paper [20]. We crop the images to center-position the object of interest and mask the background using the provided information. We report the results of this evaluation in Table 4. Evidently, we outperform PSGN and the baseline by a large margin in both Chamfer as well as EMD metrics, demonstrating the effectiveness of our approach on real data. Fig. 4 shows sample reconstructions on this dataset. Our proposed method is able to generalize well to the real dataset while both PSGN and the baseline struggle to generate meaningful reconstructions.

Category	Chamfer			EMD		
	Baseline	PSGN [8]	3D-LMNet	Baseline	PSGN [8]	3D-LMNet
airplane	3.61	3.74	3.34	7.42	6.38	4.77
bench	4.70	4.63	4.55	5.66	5.88	4.99
cabinet	7.42	6.98	6.09	9.58	6.04	6.35
car	4.67	5.20	4.55	4.74	4.87	4.10
chair	6.51	6.39	6.41	8.99	9.63	8.02
lamp	7.32	6.33	7.10	20.96	16.17	15.80
monitor	6.76	6.15	6.40	9.18	7.59	7.13
rifle	2.99	2.91	2.75	9.30	8.48	6.08
sofa	6.11	6.98	5.85	6.40	7.42	5.65
speaker	9.05	8.75	8.10	11.29	8.70	9.15
table	6.16	6.00	6.05	9.51	8.40	7.82
telephone	5.13	4.56	4.63	8.64	5.07	5.43
vessel	4.70	4.38	4.37	7.88	6.18	5.68
mean	5.78	5.62	5.40	9.20	7.75	7.00

Table 3: Single view reconstruction results on ShapeNet [8]. The metrics are computed on 1024 points after performing ICP alignment with the ground truth point cloud. All metrics are scaled by 100.

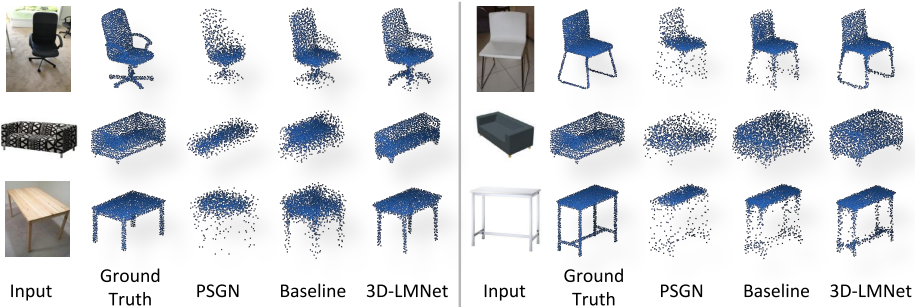


Figure 4: Qualitative results on the real-world Pix3D dataset [20]. The learnt 3D prior enables our method to generate meaningful reconstructions on real data, while both PSGN [8] and the baseline fail to generalize well to the real data distribution. PSGN predicts 1024 points, while the baseline and 3D-LMNet predict 2048 points.

Category	Chamfer			EMD		
	Baseline	PSGN [8]	3D-LMNet	Baseline	PSGN [8]	3D-LMNet
chair	7.52	8.05	7.35	11.17	12.55	9.14
sofa	8.65	8.45	8.18	8.87	9.16	7.22
table	11.23	10.82	11.20	15.71	15.16	12.73
mean	9.13	9.11	8.91	11.92	12.29	9.70

Table 4: Single view reconstruction results on the real world Pix3D dataset [20]. The metrics are computed on 1024 points after performing ICP alignment with the ground truth point cloud. All metrics are scaled by 100.

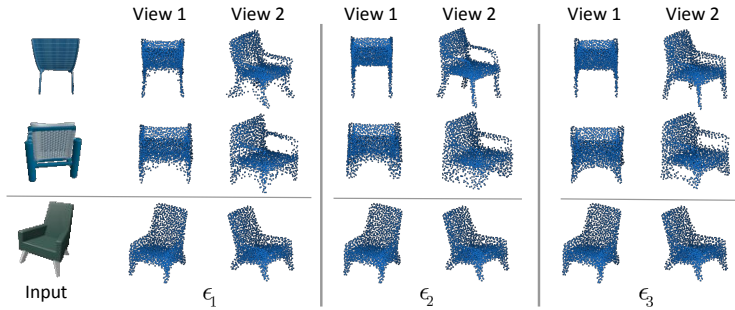


Figure 5: Qualitative results for probabilistic latent matching. Rows 1 and 2 depict multiple reconstructions for ambiguous input views, obtained by sampling ϵ from $\mathcal{N}(0, 1)$. Row 3 depicts the minimal influence ϵ has on reconstructions from informative views. Reconstruction results are shown from two different viewing angles for each ϵ so as to highlight the correspondence with the input image.

4.3 Generating multiple plausible outputs

We evaluate the probabilistic training regime (Variant II, Fig. 2c) described in Sec. 3.3 for the task of generating multiple plausible outputs for a single input image. We train E_l on objects from the chair category, and set ϕ_o and δ in Eqn. 2 to 180° and 20° respectively. For chairs, ϕ_o of 180° corresponds to a perfect back-view having maximum occlusion. For comparison, we also train a model without the diversity loss (Variant I, Fig. 2b). Quantitatively, Variant II compares favourably to Variant I in terms of both Chamfer (Variant II - 6.45 vs Variant I - 6.48) and EMD errors (Variant II - 8.04 vs Variant I - 8.1), while also effectively handling uncertainty. Qualitative results for Variant II are shown in Fig. 5. We observe that for different values of the sampling variable $\epsilon \sim \mathcal{N}(0, 1)$, we obtain semantically different reconstructions which are consistent with the input image for ambiguous views. We observe variations like presence and absence of handles, different leg structures, hollow backs, etc in the reconstructions. On the other hand, the value of ϵ has minimal influence over the reconstructions for unambiguous views.

5 Conclusion

In this paper, we highlighted the importance of learning a rich latent representation of 3D point clouds for the task of single-view 3D reconstruction. We empirically evaluated various loss formulations to effectively map to the learned latent space. We also presented a technique to tackle the inherent ambiguity in 3D shape prediction from a single image by introducing a probabilistic training scheme in the image encoder, thereby obtaining multiple plausible 3D generations from a single input image. Quantitative and qualitative evaluation on the single-image 3D reconstruction task on synthetic and real datasets show that the generated point clouds are more accurate and realistic in comparison to the current state-of-art reconstruction methods.

References

- [1] Panos Achlioptas, Olga Diamanti, Ioannis Mitliagkas, and Leonidas Guibas. Representation learning and adversarial generation of 3D point clouds. In *ICML*, 2018.
- [2] Paul J Besl and Neil D McKay. Method for registration of 3-d shapes. In *Sensor Fusion IV: Control Paradigms and Data Structures*, volume 1611, pages 586–607. International Society for Optics and Photonics, 1992.
- [3] Angel X Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, et al. Shapenet: An information-rich 3D model repository. *arXiv preprint arXiv:1512.03012*, 2015.
- [4] Christopher B Choy, Danfei Xu, JunYoung Gwak, Kevin Chen, and Silvio Savarese. 3D-r2n2: A unified approach for single and multi-view 3D object reconstruction. In *European Conference on Computer Vision*, pages 628–644. Springer, 2016.
- [5] Angela Dai, Charles Ruizhongtai Qi, and Matthias Nießner. Shape completion using 3D-encoder-predictor cnns and shape synthesis. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, volume 3, 2017.
- [6] Haowen Deng, Tolga Birdal, and Slobodan Ilic. Ppfnet: Global context aware local features for robust 3D point matching. In *CVPR*, 2018.
- [7] Carl Doersch. Tutorial on variational autoencoders. *arXiv preprint arXiv:1606.05908*, 2016.
- [8] Haoqiang Fan, Hao Su, and Leonidas Guibas. A point set generation network for 3D object reconstruction from a single image. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, volume 38, 2017.
- [9] Rohit Girdhar, David F Fouhey, Mikel Rodriguez, and Abhinav Gupta. Learning a predictable and generative vector representation for objects. In *European Conference on Computer Vision*, pages 484–499. Springer, 2016.
- [10] Christian Häne, Shubham Tulsiani, and Jitendra Malik. Hierarchical surface prediction for 3D object reconstruction. In *3DV*. 2017.
- [11] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- [12] Roman Klokov and Victor Lempitsky. Escape from cells: Deep kd-networks for the recognition of 3D point cloud models. In *2017 IEEE International Conference on Computer Vision (ICCV)*, pages 863–872. IEEE, 2017.
- [13] Andrey Kurenkov, Jingwei Ji, Animesh Garg, Viraj Mehta, JunYoung Gwak, Christopher Choy, and Silvio Savarese. Deformnet: Free-form deformation network for 3d shape reconstruction from a single image. In *WACV*, 2018.
- [14] Chen-Hsuan Lin, Chen Kong, and Simon Lucey. Learning efficient point cloud generation for dense 3D object reconstruction. In *AAAI Conference on Artificial Intelligence (AAAI)*, 2018.

- [15] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. Pointnet: Deep learning on point sets for 3D classification and segmentation. *Proc. Computer Vision and Pattern Recognition (CVPR), IEEE*, 1(2):4, 2017.
- [16] Charles Ruizhongtai Qi, Li Yi, Hao Su, and Leonidas J Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. In *Advances in Neural Information Processing Systems*, pages 5105–5114, 2017.
- [17] Danilo Jimenez Rezende, SM Ali Eslami, Shakir Mohamed, Peter Battaglia, Max Jaderberg, and Nicolas Heess. Unsupervised learning of 3D structure from images. In *Advances In Neural Information Processing Systems*, pages 4996–5004, 2016.
- [18] Gernot Riegler, Ali Osman Ulusoy, and Andreas Geiger. Octnet: Learning deep 3D representations at high resolutions. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [19] Abhishek Sharma, Oliver Grau, and Mario Fritz. Vconv-dae: Deep volumetric shape learning without object labels. In *European Conference on Computer Vision*, pages 236–250. Springer, 2016.
- [20] Xingyuan Sun, Jiajun Wu, Xiuming Zhang, Zhoutong Zhang, Chengkai Zhang, Tianfan Xue, Joshua B Tenenbaum, and William T Freeman. Pix3D: Dataset and methods for single-image 3D shape modeling. In *CVPR*, 2018.
- [21] Maxim Tatarchenko, Alexey Dosovitskiy, and Thomas Brox. Octree generating networks: Efficient convolutional architectures for high-resolution 3D outputs. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2088–2096, 2017.
- [22] Shubham Tulsiani, Tinghui Zhou, Alexei A Efros, and Jitendra Malik. Multi-view supervision for single-view reconstruction via differentiable ray consistency. In *CVPR*, volume 1, page 3, 2017.
- [23] Jiajun Wu, Chengkai Zhang, Tianfan Xue, Bill Freeman, and Josh Tenenbaum. Learning a probabilistic latent space of object shapes via 3D generative-adversarial modeling. In *Advances in Neural Information Processing Systems*, pages 82–90, 2016.
- [24] Zhirong Wu, Shuran Song, Aditya Khosla, Fisher Yu, Linguang Zhang, Xiaoou Tang, and Jianxiong Xiao. 3D shapenets: A deep representation for volumetric shapes. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1912–1920, 2015.
- [25] Xinchun Yan, Jimei Yang, Ersin Yumer, Yijie Guo, and Honglak Lee. Perspective transformer nets: Learning single-view 3D object reconstruction without 3D supervision. In *Advances in Neural Information Processing Systems*, pages 1696–1704, 2016.
- [26] Lequan Yu, Xianzhi Li, Chi-Wing Fu, Daniel Cohen-Or, and Pheng-Ann Heng. PU-Net: Point cloud upsampling network. In *CVPR*, 2018.
- [27] M Ersin Yumer and Niloy J Mitra. Learning semantic deformation flows with 3D convolutional networks. In *European Conference on Computer Vision*, pages 294–311. Springer, 2016.

Supplementary Material

1 Training Dataset Details

We train all our networks on synthetic models from the ShapeNet [1] dataset. We use the same 80%-20% train/test split provided by [1] consisting of models from 13 different categories, so as to be comparable with the previous works. We use the input images provided by [1], where each model is pre-rendered from 24 different azimuth angles. We crop the images to 128×128 resolution before passing it through our network. For generating the ground truth point cloud, we uniformly sample 16,384 points on the mesh surface using farthest point sampling.

2 Network Architectures

We provide network architecture details for the point cloud and image encoders and the common decoder in Tables 1, 2 and 3. It should be noted that 3D-LMNet has a total of 22.7M parameters, while PSGN [1] has nearly double the number with 42.9M parameters.

S.No.	Layer	Filter Size	Output Size	Params
1	conv	1x1	2048x64	0.4K
2	conv	1x1	2048x128	8.6K
3	conv	1x1	2048x128	16.8K
4	conv	1x1	2048x256	33.5K
5	conv	1x1	2048x512	132.6K
6	maxpool	-	512	0

Table 1: Point Cloud Encoder Architecture

S.No.	Layer	Filter Size/ Stride	Output Size
1	conv	3x3/1	64x64x32
2	conv	3x3/1	64x64x32
3	conv	3x3/2	32x32x64
4	conv	3x3/1	32x32x64
5	conv	3x3/1	32x32x64
6	conv	3x3/2	16x16x128
7	conv	3x3/1	16x16x128
8	conv	3x3/1	16x16x128
9	conv	3x3/2	8x8x256
10	conv	3x3/1	8x8x256
11	conv	3x3/1	8x8x256
16	conv	5x5/2	4x4x512
17	linear	-	128

Table 2: Image Encoder Architecture

S.No.	Layer	Output Size
1	linear	256
2	linear	256
3	linear	1024*3

Table 3: Decoder Architecture

Category	Chamfer			EMD		
	3D-LMNet	3D-LMNet	3D-LMNet	3D-LMNet	3D-LMNet	3D-LMNet
	Chamfer	\mathcal{L}_2	\mathcal{L}_1	Chamfer	\mathcal{L}_2	\mathcal{L}_1
airplane	4.47	3.39	3.34	7.35	4.81	4.77
bench	5.03	4.74	4.55	5.38	5.17	4.99
cabinet	6.76	6.26	6.09	7.03	6.73	6.35
car	4.70	4.61	4.55	4.31	4.20	4.10
chair	6.72	6.54	6.41	8.16	8.11	8.02
lamp	8.31	7.28	7.10	17.21	16.03	15.80
monitor	6.96	6.65	6.40	7.66	7.53	7.13
rifle	3.03	2.79	2.75	6.67	6.06	6.08
sofa	6.20	6.00	5.85	5.97	5.80	5.65
speaker	8.77	8.33	8.10	9.20	9.61	9.15
table	6.59	6.16	6.05	8.34	7.95	7.82
telephone	5.62	4.87	4.63	7.50	5.79	5.43
vessel	4.76	4.45	4.37	6.92	5.84	5.68
mean	5.99	5.54	5.40	7.82	7.20	7.00

Table 4: Category-wise 3D reconstruction metrics for different latent matching variants of 3D-LMNet on the ShapeNet dataset [9]. All metrics are scaled by 100.

3 Quantitative Comparison of 3D-LMNet Variants on ShapeNet

We report the category-wise Chamfer and EMD error metrics for all our latent matching variants on the validation split provided by [9] for the ShapeNet dataset [9] in Table 4. Our latent matching approaches (3D-LMNet- \mathcal{L}_1 and \mathcal{L}_2) significantly outperform the network trained directly with Chamfer loss (3D-LMNet-Chamfer). 3D-LMNet- \mathcal{L}_1 is better in all categories in Chamfer scores, and all but one category in terms of EMD scores.

4 Reconstructions on ShapeNet

Qualitative comparison with state-of-art and baseline for single-view reconstruction on ShapeNet validation set are provided in Figs. 1 and 2. Note that the samples are randomly selected.

5 Reconstructions on Pix3D

Qualitative comparison with state-of-art and baseline for single-view reconstruction on the real-world Pix3D dataset are shown in Fig. 3. Note that the samples are randomly selected.

6 Generating Multiple Plausible Outputs

We provide more examples for the probabilistic latent matching scheme explained in the paper in Fig. 4. We notice variations in legs, handles and back of the chair models.

7 Auto-Encoder Results

7.1 Reconstructions

3D point cloud reconstruction results are shown in Fig. 5. The reconstructions are very similar to the ground truth point clouds in appearance and spread.

7.2 Latent Space Interpolations

We analyze the quality of the learnt latent space of the auto-encoder by manipulating the latent vector z , and visually observing the generated reconstructions. Fig. 6 shows the resulting reconstructions as we linearly interpolate between two different models in the test set. We find that the interpolations are smooth and the intermediate reconstructions form valid models even in the cross-category setting.

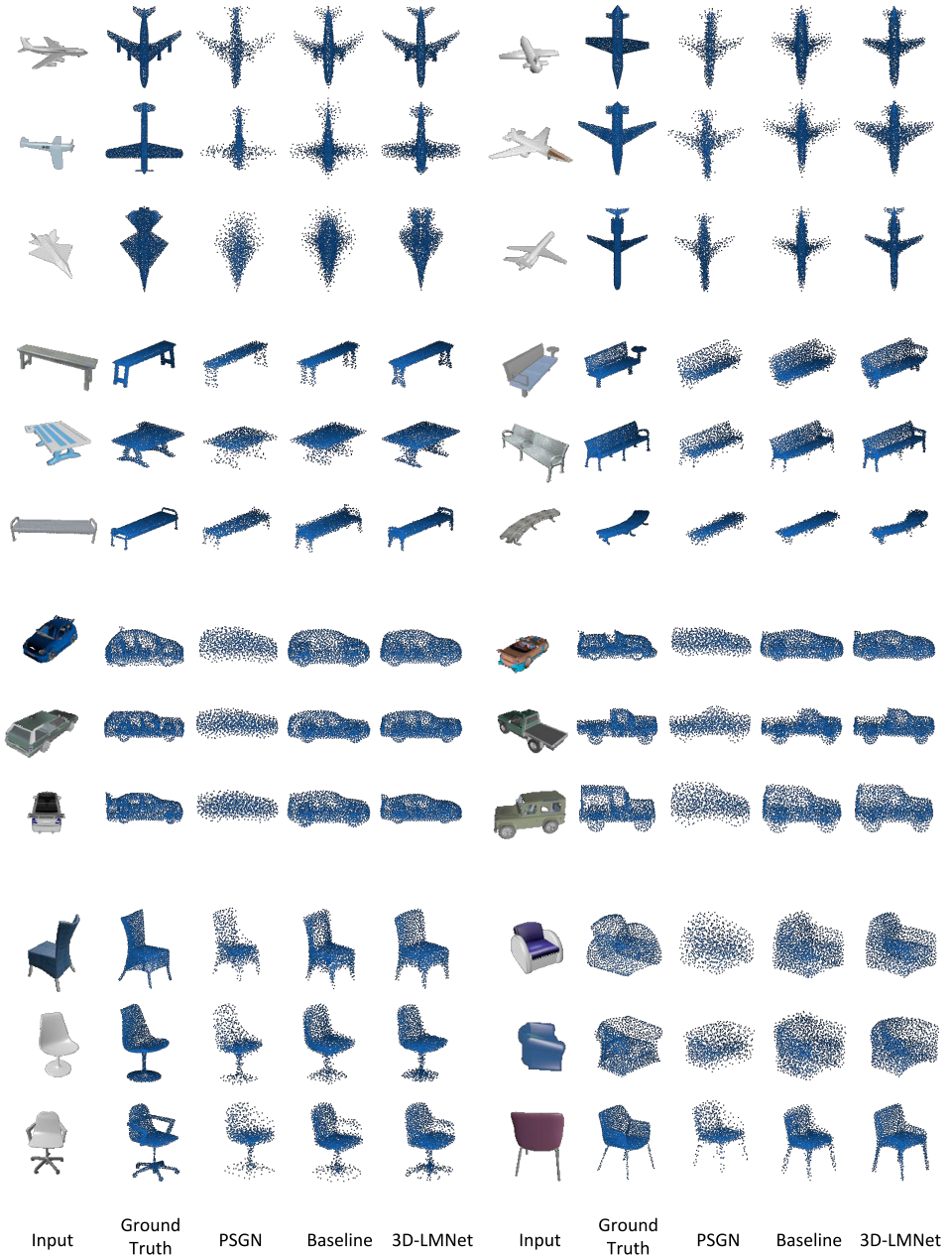


Figure 1: Reconstructions on ShapeNet. 3D reconstructions on randomly sampled input images from the validation set of ShapeNet. Note that although the baseline reconstructions for cars obtain a good shape, the points are unevenly distributed which results in high EMD error metrics (main text Table 4). On the other hand, 3D-LMNet reconstructions are well distributed and obtain lower EMD error metrics. Results best viewed zoomed.

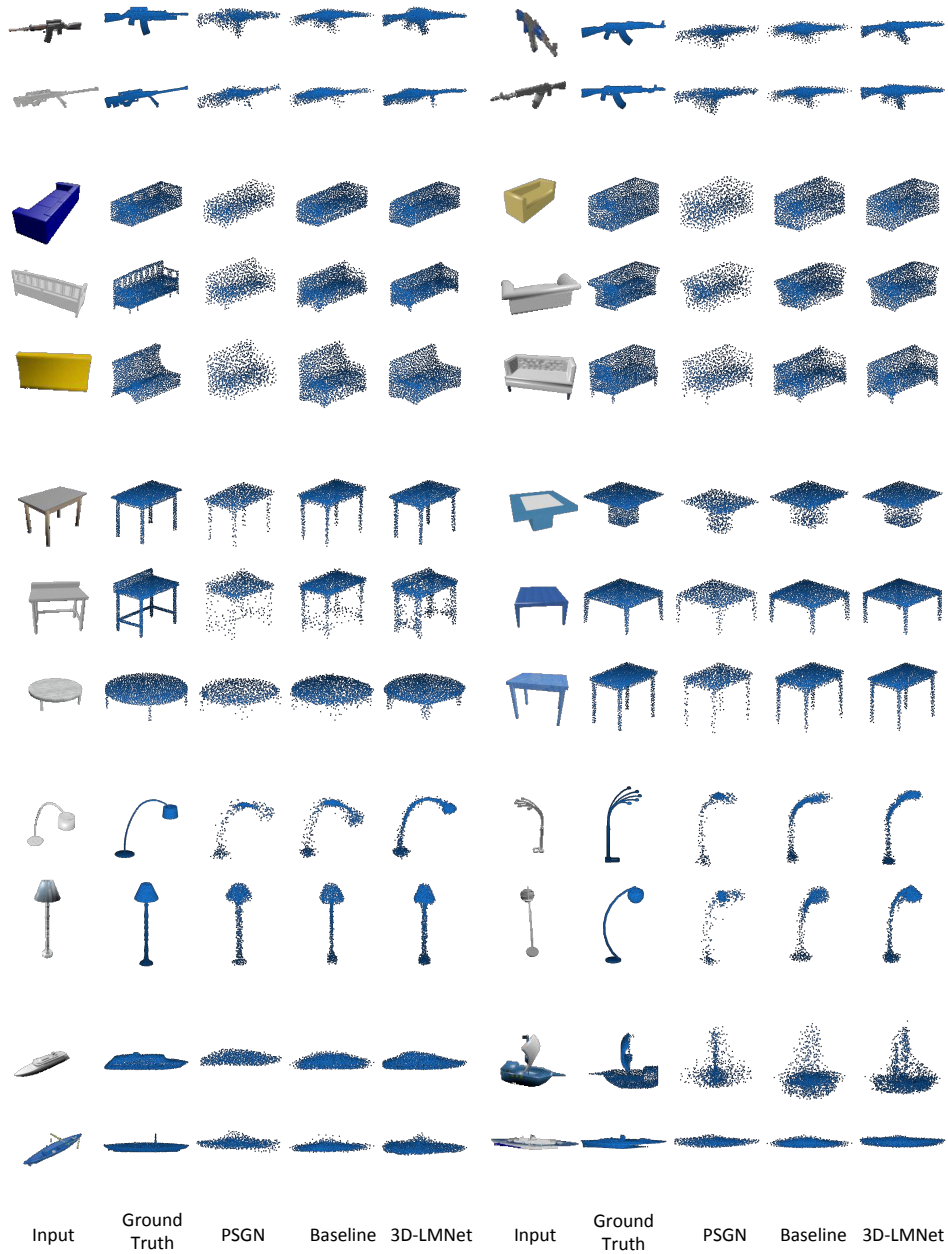


Figure 2: Reconstructions on ShapeNet. 3D reconstructions on randomly sampled input images from the validation set of ShapeNet. Results best viewed zoomed.



Figure 3: Reconstructions on Pix3d. 3D reconstructions on randomly sampled input images from Pix3D. Results best viewed zoomed.

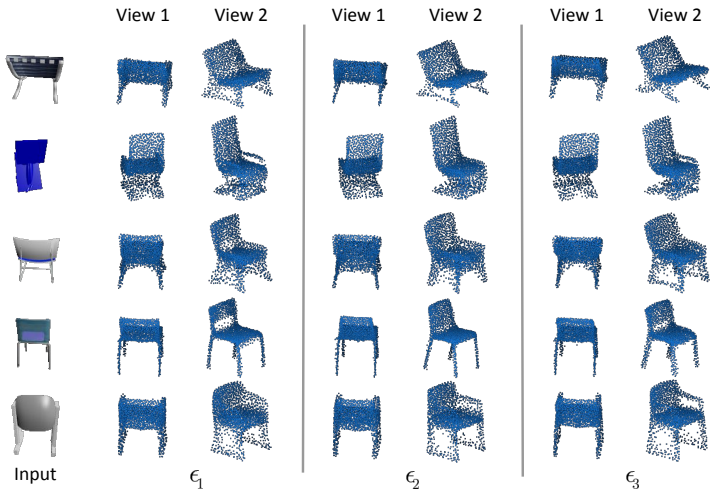


Figure 4: Qualitative results for probabilistic latent matching. Multiple reconstructions for ambiguous input views are obtained by sampling ϵ . Reconstruction results are shown from two different viewing angles for each ϵ so as to highlight the correspondence with the input image.

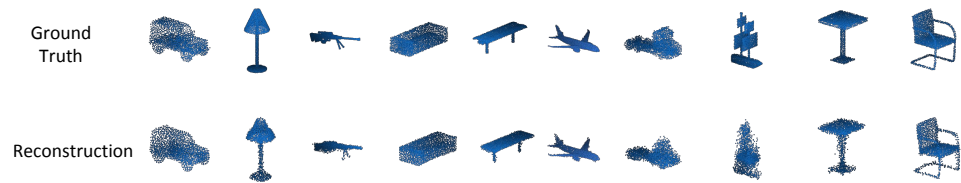


Figure 5: Auto-Encoder Reconstructions

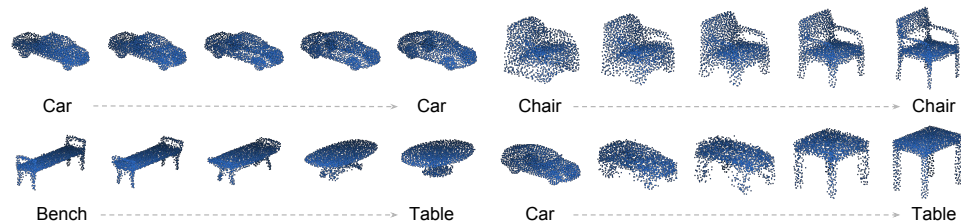


Figure 6: Auto-Encoder Interpolations