

Semantic Correspondence via 2D-3D-2D Cycle

Yang You, Chengkun Li, Yujing Lou, Zhoujun Cheng,
Lizhuang Ma, Cewu Lu, Weiming Wang*

Shanghai Jiao Tong University, China
{qq456cvb,sjtu1ck,louyujing,blankcheng}@sjtu.edu.cn,
{ma-lz,lucewu,wangweiming}@sjtu.edu.cn

Abstract. Visual semantic correspondence is an important topic in computer vision and could help machine understand objects in our daily life. However, most previous methods directly train on correspondences in 2D images, which is end-to-end but loses plenty of information in 3D spaces. In this paper, we propose a new method on predicting semantic correspondences by leveraging it to 3D domain and then project corresponding 3D models back to 2D domain, with their semantic labels. Our method leverages the advantages in 3D vision and can explicitly reason about objects self-occlusion and visibility. We show that our method gives comparative and even superior results on standard semantic benchmarks. We also conduct thorough and detailed experiments to analyze our network components. The code and experiments are publicly available at <https://github.com/qq456cvb/SemanticTransfer>.

Keywords: 3D reconstruction, deep learning, semantic transfer, differentiable renderer

1 Introduction

Semantic correspondence for general objects is an important research area for machine vision. Understanding different objects of the same category is still a challenge topic. There are quite a few methods on solving this problem in 2D image domain. [5,15,20,17,28] propose to matching local regions between pairs of images while [42,43,47,14] consider it as a global image alignment problem. However, these works all investigate 2D image features and very few works focus on the internal 3D structures of images to be matched. We argue that by explicitly exploiting 3D structures of objects, one can easily infer the self-occlusion and spatial relationships. This idea is explored in some recent works [23,63], which come up with a 3D model as an intermediate medium. However, they assume there exists a template model for all images, which does not hold in most cases.

To solve these problems, we propose a novel semantic transfer method that aims to predict 3D structures from a single RGB image and then project 3D semantic labels back onto 2D image planes. 2D to 3D shape prediction is inspired by Wu et al. [57]. For 3D-2D projection, we estimate viewpoints directly

* Weiming Wang is the corresponding author.

from 2D images and then leverage a 3D semantic prediction model trained on KeypointNet [62] to give 3D semantic labels together with its 2D projections. Viewpoints are further fine-tuned with differentiable renderers. The main advantages of this method lies in two aspects. 1) the number of training data required is reduced drastically. Previous 2D image transfer networks require numerous images for a class of object in order to extract robust semantic features. These images view objects of different shapes from different angles. On the contrary, if we could infer 3D structures from 2D images, then all we need is to utilize labels on existing 3D models and then project them onto 2D image planes. One may consider 3D structures inference as a data-heavy task but virtual 2D images can be generated from 3D models on the fly, as done by Wu et al. [57]. 2) Visibility reasoning is made explicit. When we project semantic labels onto 2D image planes, points on the back is naturally culled. On 2D image domains, visibility is implicitly made by 2D CNNs, making it hard to interpret. As shown in Figure 1, for direct 2D-2D methods, the generated 2D warping from source to target does not account for self-occlusion and may be erroneous. However, semantic transfer would be much easier if we first estimate their corresponding 3D models and camera poses.

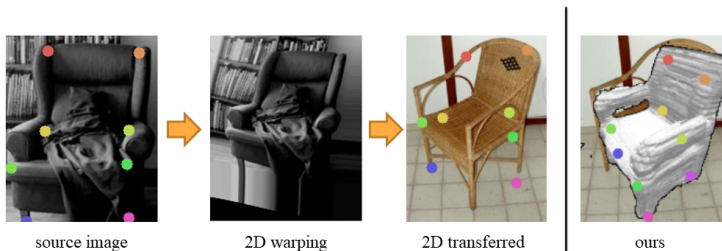


Fig. 1. Direct 2D-2D correspondence vs. ours 2D-3D-2D pipeline. Left: direct 2D semantic warping. It is erroneous to directly warp one image to another in 2D domain due to the existence of self-occlusion and viewpoint variations. Right: our method estimates the corresponding 3D model and projects predicted 3D semantic labels back onto 2D images.

Although this idea is appealing, there are still a lot of challenges in this 2D-3D-2D cycle. Compared with directly learning correspondence maps from 2D images, our pipeline is more involved. Each stage would incur some error and the final result would get biased through error accumulation. To this end, we propose a camera pose estimation module fine-tuned by differentiable renderer. The final result is competitive and even outperforms state-of-the-art methods on some benchmarks. In addition, we conduct comprehensive and detailed experiments to figure out the effectiveness of each stage in the proposed 2D-3D-2D cycle. We hope this analysis could help future researches to further improve on 2D-3D, 3D-2D or 2D-3D-2D predictions.

Our main contributions are listed below:

- We propose a novel 2D-3D-2D pipeline to solve 2D semantic correspondence problem by leveraging it to the 3D domain.
- Our proposed method sets a new state of the art on several semantic correspondence benchmarks.
- We conduct detailed and comprehensive experiments to decompose each stage/parts effect on the accuracy of final results.
- We will make all our code with detailed experiments publicly available.

2 Related Work

2.1 2D Semantic Correspondence

Image semantic correspondence has a long history which dates back to optical flow [17], multi-stereo [36]. Recently, some local descriptor based methods like proposal flow [14] and SIFT flow [28] are explored to find dense correspondences across different objects. With the advance of deep learning, neural features [16,27,22] are broadly used as they are more robust and generalizable. Methods like A2Net [47], NC-Net [44] and HPF [32] view semantic correspondence as a matching problem in high-dimensional feature images. In addition, [8,46] leverage an unsupervised methods to learn consistent dense embeddings with SLAM across different objects.

2.2 3D Semantic Correspondence

[1,2] are the pioneers on detecting 3D semantic correspondence between human bodies and faces. Recently, [12,45,11] propose unsupervised methods on learning dense correspondences between humans and animals. With the help of recent large scale model dataset such as ShapeNet [3] and PartNet [34], finding semantic correspondences on general objects become possible. Deep functional dictionaries [50] and SyncSpecCNN [60] all learn a set of synchronized base functions in order to obtain dense correspondence from functional maps. In addition to ShapeNet, [39,21,61,62] provide additional keypoint or correspondence annotations for object semantic understandings.

Perhaps, CSM [23] and Zhou et al. [63] are the closest to this paper. However, they assume that for all images, there is a template 3D model that fits well, making them not directly applicable to categories where the shapes across instances differ significantly in topology or undergo large articulation. Besides, they implicitly infer 3D models by generating a 2D-3D pixel maps while we explicitly predict each image’s corresponding 3D shape.

2.3 Single View Shape Reconstruction

Recently, many works have been introduced on single view shape reconstruction. For supervised methods where a ground-truth model is available, PSGN [7]

and pseudo-renderer [26] reconstruct point clouds from single-view RGB images. Front2back [59] predicts per-pixel depth, which is then converted into a point cloud. [55,4,56] predict voxel grids with a relatively small resolution while some others like [38,31,29] reconstruct implicit surface functions, where resolutions are not limited compared to voxels. In addition, there are also plenty of researches [10,54] focused on triangle mesh reconstruction, which is constrained by mesh topology. Pan et al. [37] tries to modify the mesh topology during reconstruction. What’s more, some other directions like reconstructing images as geometric primitive collections [9,52] and complex octree structures [41,51] are also explored.

For unsupervised single view shape prediction, [48,35,40,6] utilize only 2D image annotations, together with a multi-view consistency prior, to reconstruct the implicit 3D models. Other works like [53,25] focus on a large collection of images in the wild and reconstruct a model for each distinguished image.

2.4 Differentiable Renderer

Differentiable renderer is an emerging topic in recent years, we see that [48,35] all utilize differentiable projections to learn a 3D shape from its 2D image projections. Neural Mesh Renderer [19] first brings this idea to mesh rasterization rendering and Li et al. [24] comes up with differentiable monte-carlo ray tracing. DiffSDF [18] renders implicit surfaces defined by signed distance function in a differentiable way.

3 Methodology

3.1 Overview

Our method includes four essential parts: (a) 2D-3D shape prediction, (b) viewpoint estimation from RGB images, (c) keypoint database with semantic embeddings, (d) 3D-2D projections of semantic points. Full pipeline is illustrated in Figure 2.

For 2D-3D shape prediction, we utilize a similar structure with ShapeHD [57], which first estimate silhouettes, normals and depths from 2D images and then predict 3D shapes using 3D convolutions. This pattern can be summarized as 2D-2.5D-3D and is first proposed by MarrNet[55]. For viewpoint estimation, azimuth and elevation are classified into discrete bins using ResNet architecture. For 3D keypoint transfer, we utilize KeypointNet [62] annotations with nearest neighbor search. For 3D-2D semantic projection, 3D voxel predictions are converted into meshes with marching cube algorithm [30]. Then, these meshes can be projected back onto 2D image films provided viewpoint estimations, finetuned by differentiable renderer. Note that we assume that objects are centered in the image and do not get occluded or cut by other objects. Clutter occlusions and incompleteness are out of the scope of this paper.

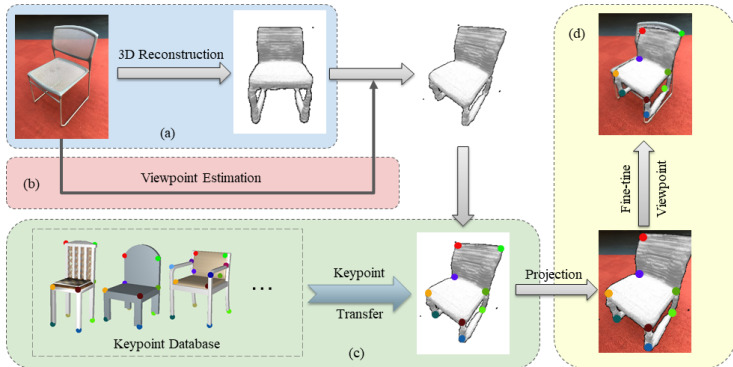


Fig. 2. Our full pipeline. (a) 3D models are reconstructed from single RGB images. (b) Viewpoints are also estimated from RGB images. (c) We obtain a keypoint descriptor database by training on existing 3D keypoint datasets, and then transfer these keypoints with nearest neighbor search. (d) Given viewpoints, 3D models and transferred keypoints, we project them onto the original image plane.

3.2 Single View 3D Reconstruction

There are a number of works focusing on single view 3D reconstruction, such as MarrNet [55], ShapeHD [57], Mesh R-CNN [10]. We utilize a similar architecture with ShapeHD, considering that it could penalize those unrealistic 3D shapes. To make it complete, here we briefly show the components that are used in ShapeHD. ShapeHD is inspired by MarrNet. It has a 2.5D sketch estimator, which is an encoder-decoder that predicts the object’s depth, surface normals and silhouette from an RGB image. Followed is a 3D estimator which also has an encoder-decoder structure. It predicts a 3D shape of the object in the canonical view from 2.5D sketches. In addition, the author introduced a deep naturalness regularizer that penalizes unrealistic shapes prediction. The regularizer is implemented by a 3D generative adversarial network and the discriminator is then used to calculate the naturalness score. The architecture for this module is shown in Figure 3.

3.3 Viewpoint Estimation from RGB Images

Given predicted 3D shapes, it is necessary to estimate the viewpoint in order to project it back onto the image plane. To do so, we design a network that predict viewpoints directly from RGB images. Note that this is different from Pix3D [49] where viewpoints are estimated from 2.5D sketches. We argue that error would accumulate if the previous 2.5D sketch prediction is inaccurate. Direct estimation reduces the number of passed stage from two to one, which help improve the accuracy. This is also verified in our experiments.

We treat view estimation as a classification problem, where azimuth is divided into 24 bins the elevation is divided into 12 bins. Circularity in azimuth is dealt

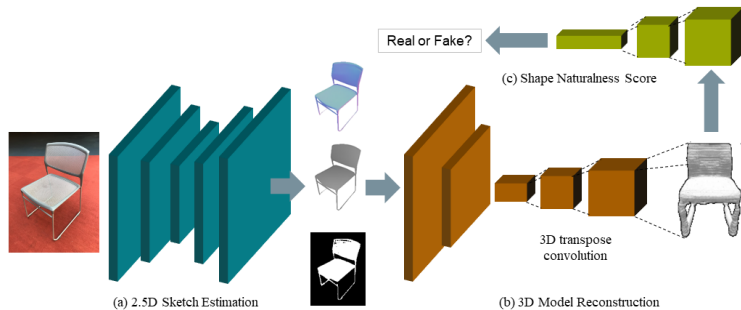


Fig. 3. Single view 3D model reconstruction. (a) Firstly, 2.5D sketches including normals, silhouettes and depths are estimated. (b) Then, 3D transpose convolution is used to recover object voxels. (c) In addition, a shape naturalness score is proposed to ensure that generated shapes are not diverged from real shapes.

carefully with an additional circular bin. The architecture for this module is demonstrated in Figure 4.

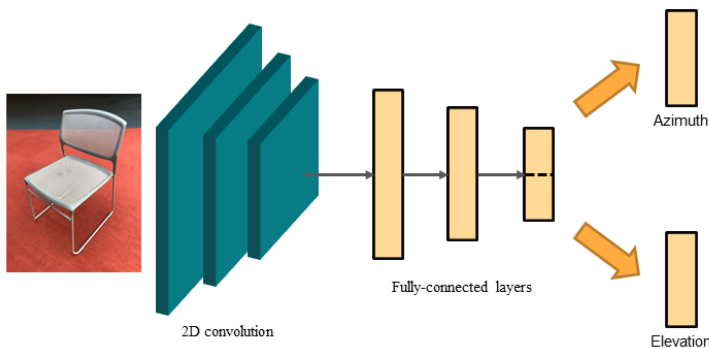


Fig. 4. Viewpoint Estimation Module. Viewpoint is estimated from 2D CNN followed by several fully connected layers. Then, KL divergence loss is employed.

3.4 3D Semantic Prediction

Predicting semantic labels on 3D models is pretty challenging in this 2D-3D-2D loop. Firstly, the predicted 3D shape from previous stage is not perfect and may be corrupted. Secondly, directly training on 3D models may be prohibitive as current semantic image datasets usually do not come up with the corresponding 3D models. Even for datasets with 3D models like PASCAL-3D [58], since the number of models is relatively small, overfitting is highly suspected.

Therefore, we resort to existing large-scale 3D keypoint datasets to train a semantic prediction network. KeypointNet[62] contains millions of keypoint annotations from ShapeNet models. By training on this dataset, one could obtain a semantic embedding for each keypoint in the dataset and then transfer it to the predicted 3D models with a nearest neighbor search. In other words, we train a semantic prediction network on a 3D object database and then generalize it to our predicted 3D shapes. To account for corruption, we augment our dataset with random Gaussian noises near the object surface. This idea is illustrated in Figure 5.

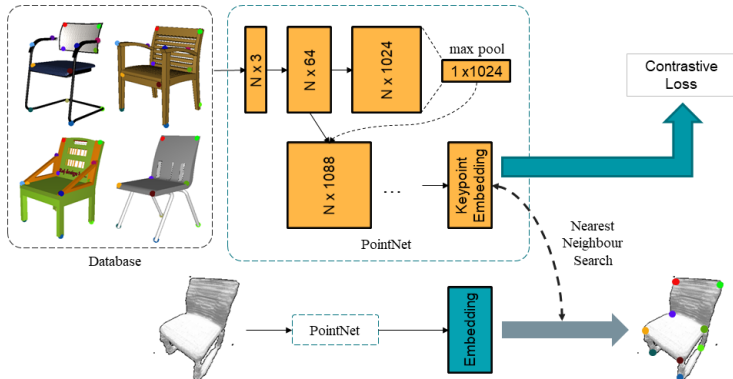


Fig. 5. Semantic Keypoint Transfer. The database has a large collection of models (may not necessarily contain the model to be evaluated). We extract their keypoint embeddings using PointNet trained with contrastive loss. For the input model, its dense embeddings are extracted with the same pretrained PointNet. Afterwards, keypoint locations are identified by a nearest neighbor search.

3.5 Differentiable Rendering in 2D Projection

As a final step, we are now ready to project our predicted 3D shapes together with inferred semantic points back onto 2D image planes. This step, although the last but not the least, is important as errors are accumulated all the way through previous stages. To have a chance correcting previous predictions, we threshold the voxels with marching cube algorithm, fine-tuning the viewpoint with the help of Neural Mesh Renderer [19]. Specifically, denote the ground-truth silhouette image as S , predicted 3D model as M , our fine-tuned viewpoint is V :

$$V^* = \arg \min_V \sum_{p \in \Omega} (S[p] - \text{Proj}(M, V)[p]), \quad (1)$$

where Ω is the set of all 2D image coordinates, $\text{Proj}(M, V)$ is the projected image given 3D model M under viewpoint V . We run several gradient descent steps in order to find the best viewpoint.

To summarize, we first predict 3D shapes and viewpoints from single view RGB images; then 3D semantic keypoints or any other semantic information is transferred from an existing 3D model database to the predicted 3D shapes; finally, the predicted 3D shapes together with their semantics are projected onto 2D image planes, with viewpoints fine-tuned.

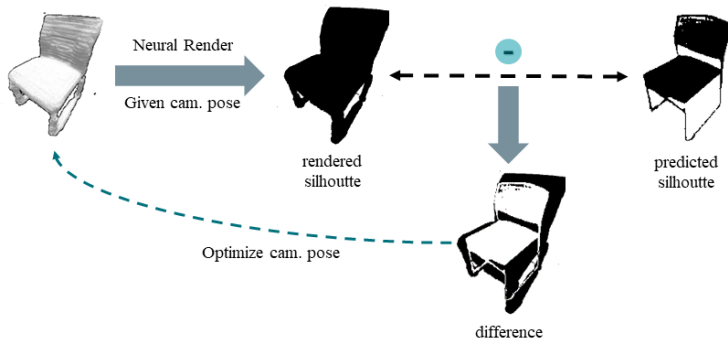


Fig. 6. Differentiable Rendering. Given current camera pose/viewpoints, we back-propagate through neural mesh renderer to optimize its pose by comparing rendered silhouette and predicted silhouette.

4 Experiments

Our experiments are divided into three parts. The first part is the comparison of our proposed method with current state-of-the-art methods. The second part is some ablation studies on our proposed viewpoint estimation and fine-tuning modules. The third part is a detailed and thorough investigation of each component/stage’s influence on final results in our full pipeline. We hope this kind of detailed analysis could help following researchers to have a better understanding on each individual components in the 2D-3D-2D loop. Note this analysis also covers previous 2D-2.5D-3D reconstruction pipeline and can be used to further improve single-view 3D reconstruction.

Datasets We use ShapeNet Synthetic rendered images [57] for training and Pix3D [49], PF-PASCAL[13], SPair-71k [33] images for evaluation (except for Section 4.4). Here, chairs are chosen for better illustrations. More results on other classes can be found in the supplementary material. Differences among these four datasets are shown in Figure 7. ShapeNet Synthetic could provide tons of training data though it is fake and synthetically rendered. Pix3D is a

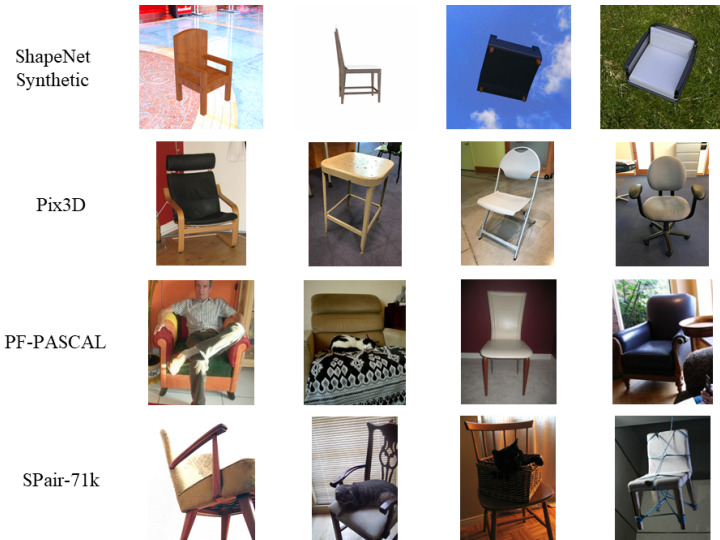


Fig. 7. Dataset Visualization (chair) on ShapeNet, Pix3D, PF-PASCAL and SPair-71k. From up to bottom: difficulties from easy to hard.

real dataset without much clutter occlusions and objects are well centered in images, which is relatively clean. PF-PASCAL and SPair-71k provide more extreme occlusions/cutoff/scale variations. Though clutter occlusions and incompletions are not the focus of this paper, our method still gives a competitive score on these datasets compared with state-of-the-art.

Metric We use a common evaluation metric of percentage of correct keypoints (PCK), which counts the average number of correctly predicted keypoints given a tolerance threshold. Given predicted keypoint \mathbf{k}_{pr} and groundtruth keypoint \mathbf{k}_{gt} , the prediction is considered correct if Euclidean distance between them is smaller than a given threshold. The correctness c of each keypoint can be expressed as

$$c = \begin{cases} 1 & \text{if } d(\mathbf{k}_{pr}, \mathbf{k}_{gt}) \leq \alpha_{\tau} \cdot \max(w_{\tau}, h_{\tau}) \\ 0 & \text{otherwise,} \end{cases} \quad (2)$$

where w_{τ} and h_{τ} are the width and height of either an entire image or object bounding box, $\tau \in \{\text{img}, \text{bbox}\}$, and α_{τ} is a tolerance factor.

We select out those keypoints that are in the intersection of both KPNet and Pix3D/PF-PASCAL/SPair-71k for evaluation. For our method, we directly estimate keypoints for each single input image and calculate PCK for each keypoint; while for other baseline methods like Hyperpixel Flow, PCK is calculated on an image pair where keypoints in one image are reckoned as ground-truth

and keypoints in the other image are predicted by a semantic warp or transfer. Our method can be also considered as transferring keypoints from an implicit ground-truth keypoint template. PCK results are averaged over all keypoints.

All our networks are written in Pytorch. Each stage in Figure 2 is trained independently. All input images are cropped and resized to 480×480 so that the object is centered in the image.

4.1 Comparison with State-of-the-Arts

In this section, we compare our methods with several state-of-the-arts that either directly do a 2D image semantic transfer [32,47] or utilize 3D templates [23].

Our method is trained with ShapeNet Synthetic renderings while state-of-the-art methods are trained directly on real-world images. Interestingly, though our method has a domain gap when applied to real-world images, we still outperform state-of-the-art methods on Spair-71k and Pix3D. Quantitative results are shown in Table 1. On PF-PASCAL, our method is inferior due to the difficulty in handling severely occluded and incomplete objects. Qualitative results on SPair-71k/PF-PASCAL are shown in Figure 8 and 9.

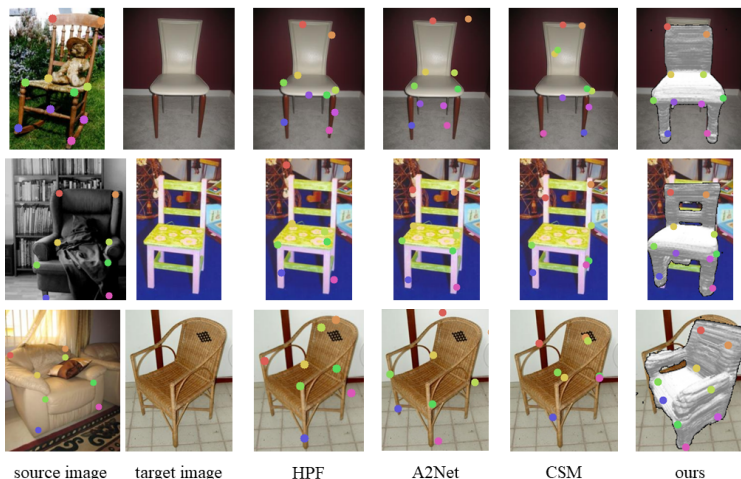


Fig. 8. Qualitative results on PF-PASCAL.

4.2 Ablation Study on Viewpoint Estimation

In this section, we explore the effect of proposed view estimation module and viewpoint fine-tuning module.

For the view estimation module, we compare with viewpoints that are predicted from estimated 2.5D sketch (w/o v.p. from RGB). Quantitative results

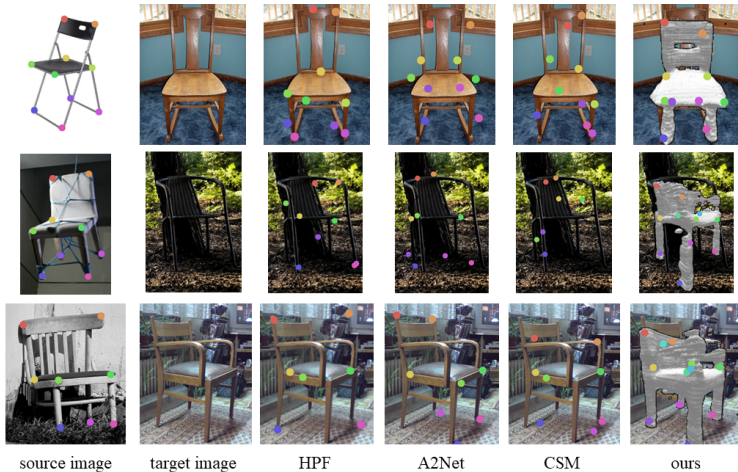


Fig. 9. Qualitative results on SPair-71k.

Models	PCK-chair ($\alpha_{img} = 0.1$)			PCK-chair ($\alpha_{bbox} = 0.1$)		
	PF-PASCAL	SPair71k	Pix3D	PF-PASCAL	SPair71k	Pix3D
ours	0.444	0.565	0.560	0.342	0.346	0.323
HPF _{res101} [32]	0.602	0.419	0.534	0.435	0.324	0.434
A2Net _{res101} [47]	0.516	0.358	0.506	0.302	0.200	0.356
CSM _{unet} [23]	0.164	0.115	0.152	0.164 ¹	0.115 ¹	0.152 ¹

Table 1. Comparison of our method with state-of-the-arts. ¹CSM crops input images with object bounding boxes, so the results for $\alpha_{img} = 0.1$ and $\alpha_{bbox} = 0.1$ are the same.

are shown in Table 2, we see that our proposed method greatly improve the accuracy of view estimation by not accumulating the error in the 2D-to-2.5D sketch prediction. From Figure 10, it can be concluded that viewpoints estimated from 2.5D sketch are much more biased and reduce the quality of final transferred keypoints.

For the viewpoint fine-tuning module, fine-tuning viewpoints improves the overall accuracy on SPair-71k and Pix3D while downgrades on PF-PASCAL when $\alpha_{img} = 0.1$, as shown in Table 2. This is due to the fact that PF-PASCAL includes more occluded chairs than SPair-71k thus is harder than the latter one. This breaks the prior of objects centered on the image thus making viewpoint fine-tuning vulnerable. Qualitative results are given in Figure 10.

4.3 Detailed Analysis of Each Stage

In this section, we investigate how each stage influences the final result, by replacing the following three components with their ground-truths: (a) 2D to 3D shape reconstruction, (b) viewpoint estimation and (c) semantic 3D model.

Models	PCK-chair ($\alpha_{img} = 0.1$)			PCK-chair ($\alpha_{bbox} = 0.1$)		
	PF-PASCAL	SPair71k	Pix3D	PF-PASCAL	SPair71k	Pix3D
ours	0.444	0.565	0.560	0.342	0.346	0.323
ours w/o v.p. from RGB	0.171	0.185	0.238	0.075	0.098	0.117
ours w/o v.p. fine-tune	0.497	0.538	0.560	0.316	0.332	0.322

Table 2. Ablation study on viewpoint estimation modules.

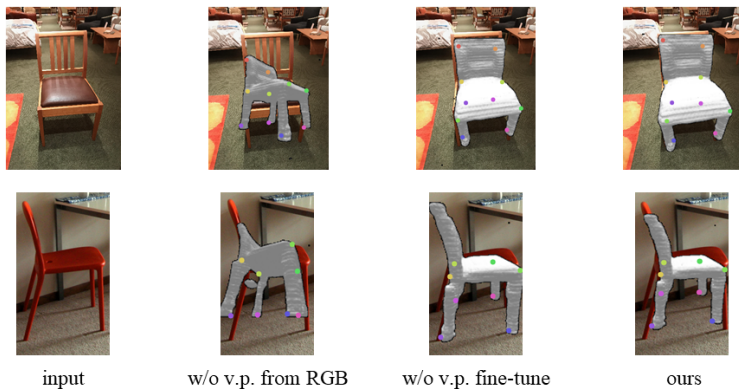


Fig. 10. Visualization on viewpoint estimation modules. From left to right: input image; viewpoint estimated from predicted 2.5D sketches instead of RGB images; viewpoint not fine-tuned; models with all modules.

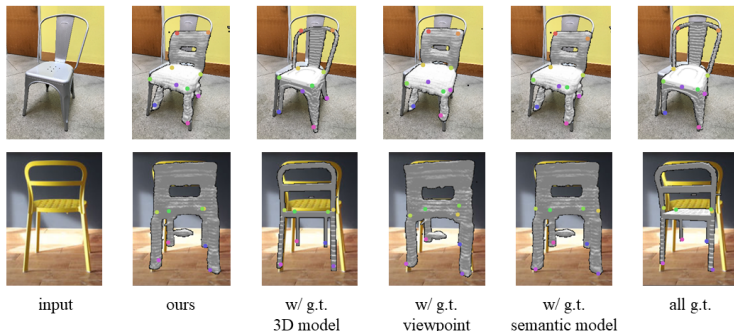


Fig. 11. Visualization of each stage's effect on Pix3D. From left to right: input image; replaced with ground-truth 3D model; replaced with ground-truth viewpoint; replaced with ground-truth 3D model in semantic transfer; all ground-truths.

Here, semantic 3D model means whether to use ground-truth 3D model of the input image when doing keypoint transfer (the database input to PointNet in Figure 5).

We start with our full pipeline and then replace each component with its ground-truth counterpart to see the accuracy improvement, respectively. We also

evaluate our method with all components’ corresponding ground-truths. Notice that although ground-truth viewpoints with azimuth and elevation are given, they are not the ground-truth 6D camera poses. Therefore, all components with ground-truth fails to achieve 100% accuracy.

Results are given in Table 3. We see that the PCK contribution of ground-truth semantic 3D model is the largest, which means that if we have the ground-truth model for computing keypoint embeddings instead of the ones in our keypoint database, we would gain about 15.5% relative improvement. The contribution of ground-truth 3D reconstruction is small, which suggests that the 2D-2.5D-3D single view reconstruction pipeline meets few difficulties when applied to real datasets. Replacing predicted viewpoints with ground-truths also gives 8.2% PCK improvement, which means that there is still some future work to do in single view camera pose estimation. More visualization results are demonstrated in Figure 11.

Models	PCK-chair ($\alpha_{img} = 0.1$)		PCK-chair ($\alpha_{bbox} = 0.1$)	
	Pix3D	ShapeNet	Pix3D	ShapeNet
ours	0.560	0.323	0.513	0.243
ours w/ 2.5D model GT	-	-	0.518	0.243
ours w/ 3D model GT	0.571	0.363	0.521	0.263
ours w/ viewpoint GT	0.606	0.351	0.545	0.262
ours w/ semantic model GT	0.647	0.434	0.631	0.331
all GT	0.741	0.560	0.723	0.419

Table 3. Detailed analysis of each stage on Pix3D and ShapeNet.

4.4 Domain Gap Exploration

Plenty of single view 3D reconstructions are done on virtually rendered datasets and validated on real-world data. This greatly reduces the need for human annotated 2D keypoint and 3D model label pairs. However, this also introduces an unavoidable gap since due to the rendering error.

In this section, we evaluate our method on virtually rendered ShapeNet, which is from the same domain of training datasets while all the evaluated images are not seen during training. This is in comparison with Section 4.1 whose evaluation is on real datasets. Note, since we are evaluating on rendered datasets, we have 2.5D ground-truth (silhouette, normal and depth), so that we add an extra experiment by replacing 2.5D predictions with ground-truth data.

Quite interestingly, results on Pix3D are better than on ShapeNet Synthetic. One reason is that in ShapeNet Synthetic dataset, camera positions are sampled uniformly from the entire unit sphere while real datasets seldom have large elevations, making generalization from virtual datasets easier. Quantitative results are given in Table 3 and visualization is shown in Figure 12.

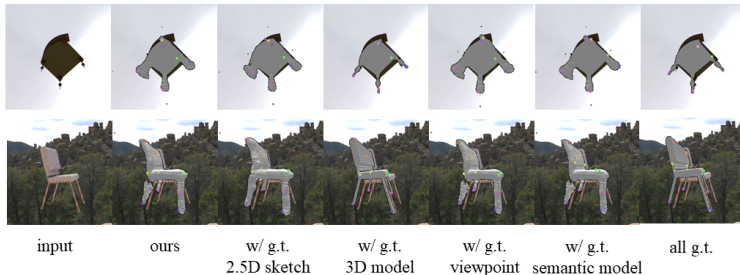


Fig. 12. Detailed analysis of each stage on ShapeNet Synthetic. From left to right: input image; replaced with ground-truth 2.5D sketch; replaced with ground-truth 3D model; replaced with ground-truth viewpoint; replaced with ground-truth 3D model in semantic transfer; all ground-truths.



Fig. 13. Failure cases of our results on SPair-71k/PF-PASCAL when there exists severe clutter occlusions.

5 Future Work

Here we show some failure cases of our method in Figure 13, where severe clutter occlusions are introduced. It would be interesting to extend our method to explicit reason about clutter occlusions and incompletions. Besides, as our current viewpoint estimation only includes two degrees of freedom, a full 6D pose estimation is possible and we leave it as a future work.

6 Conclusions

In this paper, we propose a new pipeline on predicting semantic correspondences by leveraging it to 3D domain and then project corresponding 3D models back to 2D domain, with their semantic labels. This method explicitly reasons about objects self-occlusion and visibility. We show that our method gives comparative and even superior results on standard semantic benchmarks. We also conduct thorough and detailed experiments to analyze our network components.

Supplementary

Dense Embedding Prediction

In the main document, we show that our method outperforms other methods on keypoint transfer task. Dense embeddings from single RGB images can also be obtained by propogating 3D semantic embeddings onto 2D image plane. Some qualitative results are illustrated in 14 and 15.



Fig. 14. Predicted dense embeddings on PF-PASCAL cars. Notice how the generated embeddings are consistent across different models, despite of viewpoint variations. Similar colors indicate similar embeddings.



Fig. 15. Predicted dense embeddings on PF-PASCAL aeroplanes. Notice how the generated embeddings are consistent across different models, despite of viewpoint variations. Similar color indicate similar embeddings.

Keypoint Transfer Results on Other Classes

In this section, we provide additional PCK results on both car and aeroplanes, evaluated on PF-PASCAL and SPair-71k.

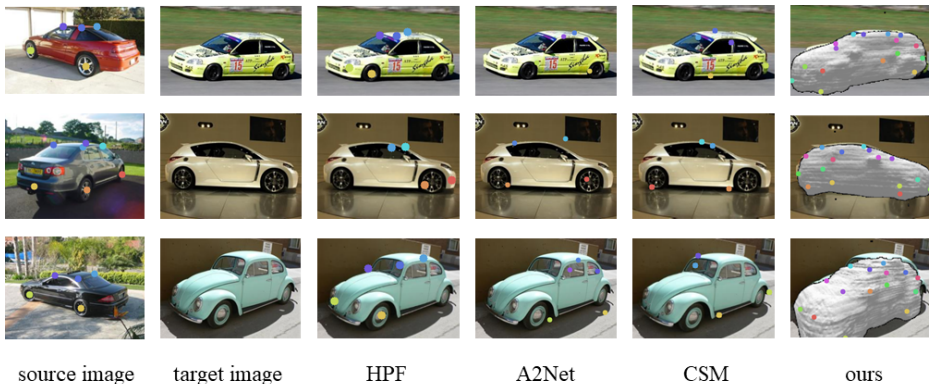


Fig. 16. Qualitative results on PF-PASCAL cars.



Fig. 17. Qualitative results on SPair-71k cars.

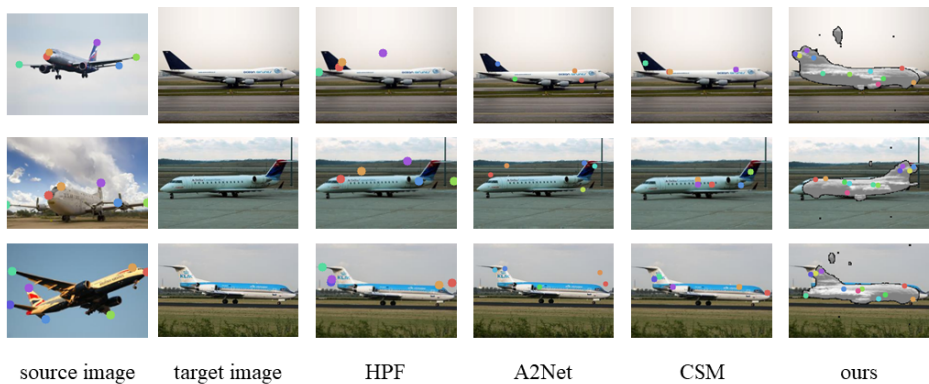


Fig. 18. Qualitative results on PF-PASCAL aeroplanes.

Models	PCK-car ($\alpha_{img} = 0.1$)		PCK-car ($\alpha_{bbox} = 0.1$)	
	PF-PASCAL	SPair71k	PF-PASCAL	SPair71k
ours	0.574	0.500	0.349	0.328
HPF _{res101}	0.533	0.364	0.437	0.276
A2Net _{res101}	0.478	0.332	0.326	0.201
CSM _{unet}	0.339	0.234	0.339	0.234

Table 4. Comparison of our method with state-of-the-arts on cars.

Models	PCK-aeroplane ($\alpha_{img} = 0.1$)		PCK-aeroplane ($\alpha_{bbox} = 0.1$)	
	PF-PASCAL	SPair71k	PF-PASCAL	SPair71k
ours	0.440	0.390	0.194	0.182
HPF _{res101}	0.401	0.280	0.294	0.212
A2Net _{res101}	0.388	0.234	0.256	0.155
CSM _{unet}	0.220	0.139	0.220	0.139

Table 5. Comparison of our method with state-of-the-arts on aeroplanes.

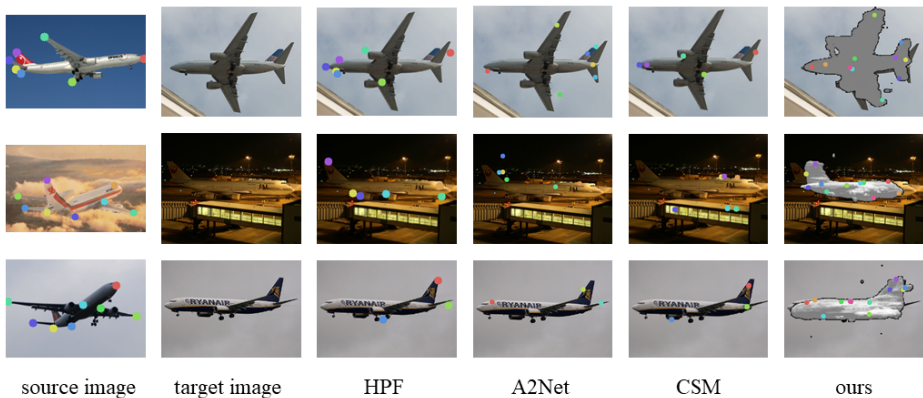


Fig. 19. Qualitative results on SPair-71k aeroplanes.

References

1. Allen, B., Curless, B., Popović, Z.: The space of human body shapes: reconstruction and parameterization from range scans. *ACM transactions on graphics (TOG)* **22**(3), 587–594 (2003)
2. Blanz, V., Vetter, T.: A morphable model for the synthesis of 3d faces. In: *Proceedings of the 26th annual conference on Computer graphics and interactive techniques*. pp. 187–194 (1999)
3. Chang, A.X., Funkhouser, T., Guibas, L., Hanrahan, P., Huang, Q., Li, Z., Savarese, S., Savva, M., Song, S., Su, H., et al.: Shapenet: An information-rich 3d model repository. *arXiv preprint arXiv:1512.03012* (2015)
4. Chen, Y., Cipolla, R.: Learning shape priors for single view reconstruction. In: *2009 IEEE 12th International Conference on Computer Vision Workshops, ICCV Workshops*. pp. 1425–1432. IEEE (2009)
5. Choy, C.B., Gwak, J., Savarese, S., Chandraker, M.: Universal correspondence network. In: *Advances in Neural Information Processing Systems*. pp. 2414–2422 (2016)
6. Eslami, S.A., Rezende, D.J., Besse, F., Viola, F., Morcos, A.S., Garnelo, M., Ruderman, A., Rusu, A.A., Danihelka, I., Gregor, K., et al.: Neural scene representation and rendering. *Science* **360**(6394), 1204–1210 (2018)
7. Fan, H., Su, H., Guibas, L.J.: A point set generation network for 3d object reconstruction from a single image. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 605–613 (2017)
8. Florence, P.R., Manuelli, L., Tedrake, R.: Dense object nets: Learning dense visual object descriptors by and for robotic manipulation. *arXiv preprint arXiv:1806.08756* (2018)
9. Gao, L., Yang, J., Wu, T., Yuan, Y.J., Fu, H., Lai, Y.K., Zhang, H.: Sdm-net: Deep generative network for structured deformable mesh. *ACM Transactions on Graphics (TOG)* **38**(6), 1–15 (2019)
10. Gkioxari, G., Malik, J., Johnson, J.: Mesh r-cnn. In: *Proceedings of the IEEE International Conference on Computer Vision*. pp. 9785–9795 (2019)
11. Groueix, T., Fisher, M., Kim, V.G., Russell, B.C., Aubry, M.: 3d-coded: 3d correspondences by deep deformation. In: *Proceedings of the European Conference on Computer Vision (ECCV)*. pp. 230–246 (2018)
12. Halimi, O., Litany, O., Rodolà, E., Bronstein, A., Kimmel, R.: Self-supervised learning of dense shape correspondence. *arXiv preprint arXiv:1812.02415* (2018)
13. Ham, B., Cho, M., Schmid, C., Ponce, J.: Proposal flow. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. pp. 3475–3484 (2016)
14. Ham, B., Cho, M., Schmid, C., Ponce, J.: Proposal flow: Semantic correspondences from object proposals. *IEEE transactions on pattern analysis and machine intelligence* **40**(7), 1711–1725 (2017)
15. Han, K., Rezende, R.S., Ham, B., Wong, K.Y.K., Cho, M., Schmid, C., Ponce, J.: Snet: Learning semantic correspondence. In: *Proceedings of the IEEE International Conference on Computer Vision*. pp. 1831–1840 (2017)
16. Hariharan, B., Arbeláez, P., Girshick, R., Malik, J.: Hypercolumns for object segmentation and fine-grained localization. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 447–456 (2015)
17. Horn, B.K., Schunck, B.G.: Determining optical flow. In: *Techniques and Applications of Image Understanding*. vol. 281, pp. 319–331. International Society for Optics and Photonics (1981)

18. Jiang, Y., Ji, D., Han, Z., Zwicker, M.: Sdfdiff: Differentiable rendering of signed distance fields for 3d shape optimization. arXiv preprint arXiv:1912.07109 (2019)
19. Kato, H., Ushiku, Y., Harada, T.: Neural 3d mesh renderer. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 3907–3916 (2018)
20. Kim, S., Min, D., Ham, B., Jeon, S., Lin, S., Sohn, K.: Fcss: Fully convolutional self-similarity for dense semantic correspondence. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 6560–6569 (2017)
21. Kim, V.G., Li, W., Mitra, N.J., Chaudhuri, S., DiVerdi, S., Funkhouser, T.: Learning part-based templates from large collections of 3d shapes. *ACM Transactions on Graphics (TOG)* **32**(4), 1–12 (2013)
22. Kong, T., Yao, A., Chen, Y., Sun, F.: Hypernet: Towards accurate region proposal generation and joint object detection. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 845–853 (2016)
23. Kulkarni, N., Gupta, A., Tulsiani, S.: Canonical surface mapping via geometric cycle consistency. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 2202–2211 (2019)
24. Li, T.M., Aittala, M., Durand, F., Lehtinen, J.: Differentiable monte carlo ray tracing through edge sampling. *ACM Transactions on Graphics (TOG)* **37**(6), 1–11 (2018)
25. Li, X., Dong, Y., Peers, P., Tong, X.: Synthesizing 3d shapes from silhouette image collections using multi-projection generative adversarial networks. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 5535–5544 (2019)
26. Lin, C.H., Kong, C., Lucey, S.: Learning efficient point cloud generation for dense 3d object reconstruction. In: Thirty-Second AAAI Conference on Artificial Intelligence (2018)
27. Lin, T.Y., Dollár, P., Girshick, R., He, K., Hariharan, B., Belongie, S.: Feature pyramid networks for object detection. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 2117–2125 (2017)
28. Liu, C., Yuen, J., Torralba, A.: Sift flow: Dense correspondence across scenes and its applications. *IEEE transactions on pattern analysis and machine intelligence* **33**(5), 978–994 (2010)
29. Liu, S., Saito, S., Chen, W., Li, H.: Learning to infer implicit surfaces without 3d supervision. In: Advances in Neural Information Processing Systems. pp. 8293–8304 (2019)
30. Lorensen, W.E., Cline, H.E.: Marching cubes: A high resolution 3d surface construction algorithm. *ACM siggraph computer graphics* **21**(4), 163–169 (1987)
31. Mescheder, L., Oechsle, M., Niemeyer, M., Nowozin, S., Geiger, A.: Occupancy networks: Learning 3d reconstruction in function space. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 4460–4470 (2019)
32. Min, J., Lee, J., Ponce, J., Cho, M.: Hyperpixel flow: Semantic correspondence with multi-layer neural features. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 3395–3404 (2019)
33. Min, J., Lee, J., Ponce, J., Cho, M.: Spair-71k: A large-scale benchmark for semantic correspondence. arXiv preprint arXiv:1908.10543 (2019)
34. Mo, K., Zhu, S., Chang, A.X., Yi, L., Tripathi, S., Guibas, L.J., Su, H.: Partnet: A large-scale benchmark for fine-grained and hierarchical part-level 3d object understanding. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 909–918 (2019)

35. Niemeyer, M., Mescheder, L., Oechsle, M., Geiger, A.: Differentiable volumetric rendering: Learning implicit 3d representations without 3d supervision. arXiv preprint arXiv:1912.07372 (2019)
36. Okutomi, M., Kanade, T.: A multiple-baseline stereo. *IEEE Transactions on pattern analysis and machine intelligence* **15**(4), 353–363 (1993)
37. Pan, J., Han, X., Chen, W., Tang, J., Jia, K.: Deep mesh reconstruction from single rgb images via topology modification networks. In: *Proceedings of the IEEE International Conference on Computer Vision*. pp. 9964–9973 (2019)
38. Park, J.J., Florence, P., Straub, J., Newcombe, R., Lovegrove, S.: DeepSDF: Learning continuous signed distance functions for shape representation. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. pp. 165–174 (2019)
39. Pavlakos, G., Zhou, X., Chan, A., Derpanis, K.G., Daniilidis, K.: 6-dof object pose from semantic keypoints. In: *2017 IEEE international conference on robotics and automation (ICRA)*. pp. 2011–2018. IEEE (2017)
40. Rematas, K., Ferrari, V.: Neural voxel renderer: Learning an accurate and controllable rendering tool. arXiv preprint arXiv:1912.04591 (2019)
41. Riegler, G., Osman Ulusoy, A., Geiger, A.: Octnet: Learning deep 3d representations at high resolutions. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. pp. 3577–3586 (2017)
42. Rocco, I., Arandjelovic, R., Sivic, J.: Convolutional neural network architecture for geometric matching. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. pp. 6148–6157 (2017)
43. Rocco, I., Arandjelović, R., Sivic, J.: End-to-end weakly-supervised semantic alignment. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. pp. 6917–6925 (2018)
44. Rocco, I., Cimpoi, M., Arandjelović, R., Torii, A., Pajdla, T., Sivic, J.: Neighbourhood consensus networks. In: *Advances in Neural Information Processing Systems*. pp. 1651–1662 (2018)
45. Roufousse, J.M., Sharma, A., Ovsjanikov, M.: Unsupervised deep learning for structured shape matching. In: *Proceedings of the IEEE International Conference on Computer Vision*. pp. 1617–1627 (2019)
46. Schmidt, T., Newcombe, R., Fox, D.: Self-supervised visual descriptor learning for dense correspondence. *IEEE Robotics and Automation Letters* **2**(2), 420–427 (2016)
47. Seo, P.H., Lee, J., Jung, D., Han, B., Cho, M.: Attentive semantic alignment with offset-aware correlation kernels. In: *European Conference on Computer Vision*. pp. 367–383. Springer (2018)
48. Sitzmann, V., Thies, J., Heide, F., Nießner, M., Wetzstein, G., Zollhofer, M.: Deepvoxels: Learning persistent 3d feature embeddings. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. pp. 2437–2446 (2019)
49. Sun, X., Wu, J., Zhang, X., Zhang, Z., Zhang, C., Xue, T., Tenenbaum, J.B., Freeman, W.T.: Pix3d: Dataset and methods for single-image 3d shape modeling. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. pp. 2974–2983 (2018)
50. Sung, M., Su, H., Yu, R., Guibas, L.J.: Deep functional dictionaries: Learning consistent semantic structures on 3d models from functions. In: *Advances in Neural Information Processing Systems*. pp. 485–495 (2018)
51. Tatarchenko, M., Dosovitskiy, A., Brox, T.: Octree generating networks: Efficient convolutional architectures for high-resolution 3d outputs. In: *Proceedings of the IEEE International Conference on Computer Vision*. pp. 2088–2096 (2017)

52. Tian, Y., Luo, A., Sun, X., Ellis, K., Freeman, W.T., Tenenbaum, J.B., Wu, J.: Learning to infer and execute 3d shape programs. arXiv preprint arXiv:1901.02875 (2019)
53. Wang, C., Lin, C.H., Lucey, S.: Deep nrsfm++: Towards 3d reconstruction in the wild. arXiv preprint arXiv:2001.10090 (2020)
54. Wen, C., Zhang, Y., Li, Z., Fu, Y.: Pixel2mesh++: Multi-view 3d mesh generation via deformation. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 1042–1051 (2019)
55. Wu, J., Wang, Y., Xue, T., Sun, X., Freeman, B., Tenenbaum, J.: Marrnet: 3d shape reconstruction via 2.5 d sketches. In: Advances in neural information processing systems. pp. 540–550 (2017)
56. Wu, J., Zhang, C., Xue, T., Freeman, B., Tenenbaum, J.: Learning a probabilistic latent space of object shapes via 3d generative-adversarial modeling. In: Advances in neural information processing systems. pp. 82–90 (2016)
57. Wu, J., Zhang, C., Zhang, X., Zhang, Z., Freeman, W.T., Tenenbaum, J.B.: Learning shape priors for single-view 3d completion and reconstruction. In: Proceedings of the European Conference on Computer Vision (ECCV). pp. 646–662 (2018)
58. Xiang, Y., Mottaghi, R., Savarese, S.: Beyond pascal: A benchmark for 3d object detection in the wild. In: IEEE Winter Conference on Applications of Computer Vision (WACV) (2014)
59. Yao, Y., Schertler, N., Rosales, E., Rhodin, H., Sigal, L., Sheffer, A.: Front2back: Single view 3d shape reconstruction via front to back prediction. arXiv preprint arXiv:1912.10589 (2019)
60. Yi, L., Su, H., Guo, X., Guibas, L.J.: Syncspecnn: Synchronized spectral cnn for 3d shape segmentation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 2282–2290 (2017)
61. You, Y., Li, C., Lou, Y., Cheng, Z., Li, L., Ma, L., Wang, W., Lu, C.: Fine-grained object semantic understanding from correspondences. arXiv preprint arXiv:1912.12577 (2019)
62. You, Y., Lou, Y., Li, C., Cheng, Z., Li, L., Ma, L., Lu, C., Wang, W.: Keypointnet: A large-scale 3d keypoint dataset aggregated from numerous human annotations (2020)
63. Zhou, T., Krahenbuhl, P., Aubry, M., Huang, Q., Efros, A.A.: Learning dense correspondence via 3d-guided cycle consistency. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 117–126 (2016)