# **Appendix: Unsupervised 3D Keypoint Discovery for Shape Control**

Tomas Jakab<sup>1,4\*</sup>, Richard Tucker<sup>4</sup>, Ameesh Makadia<sup>4</sup>, Jiajun Wu<sup>3</sup>, Noah Snavely<sup>4</sup>, Angjoo Kanazawa<sup>2,4</sup> <sup>1</sup>University of Oxford, <sup>2</sup>UC Berkeley, <sup>3</sup>Stanford University, <sup>4</sup>Google Research

This supplemental document provides an ablation study (Appendix A), extensive results (Appendices B and C), and further implementation details (Appendix D). Please also refer to our video on our project page<sup>1</sup> that demonstrates our method in action.

### A. Ablations



Figure 1: Farthest Point Keypoint regularizer ablation. We investigate the influence of the number J of sampled farthest points q used for the keypoint regularizer (Section 3.2 of the main paper) on the quality of discovered unsupervised keypoints. We show unsupervised 3D keypoints trained using two versions of regularization. First, we set the number of sampled farthest points to double of the number of keypoints (a, c). This is the setup that we use throughout the paper. Second, we set the number of sampled farthest points to the number of points in the point cloud representing the shape. This essentially results in a regularizer that is minimizing the Chamfer distance between the unsupervised keypoints and the object point cloud. Although the learned unsupervised keypoints have a good coverage (b, d) they are not as equally spaced and characteristic of the shape as (a, c).

**Varying number of regularizing points.** We examine the importance of the number of sampled farthest points q on the quality of keypoint regularization (Section 3.2 of the main paper). Figure 1 shows the effect of different numbers of sampled farthest points on the discovered keypoints. Using a high number of sampled farthest points in the regularization fails to learn keypoints that are equally spaced and characteristic of the underlying shape.

**Varying number of keypoints.** We vary the number of unsupervised keypoints discovered by our method. Figure 2 shows that our keypoints remain semantically consistent for different numbers of discovered keypoints.

<sup>\*</sup> Work done while interning at Google Research.

<sup>&</sup>lt;sup>1</sup>http://tomasjakab.github.io/KeypointDeformer

#### **B.** Shape Control via Unsupervised 3D Keypoints

We show user guided interactive shape control in our supplementary video on our project page. Figure 5 shows frames captured from user-guided interactive shape editing. Editing using our keypoints is fast and intuitive while preserving the character and details of the original shape.

#### C. Unsupervised 3D Keypoints

We show extended quantitative results for semantic part correspondence experiment and detailed correlation tables in Figure 3 for the ShapeNet Car category. Figures 6 to 10 show extensive *randomly* sampled qualitative test results for our unsupervised 3D keypoints.

#### **D.** Implementation Details

Our model assumes that the shapes are aligned (in the same orientation). The initial cage is a 42-vertex icosphere. We limit the influence matrix W to influence at most M nearest cages vertices (Section 3.1) per each keypoint, with  $M = \lfloor C/K \rfloor$ , where C is the number of cage vertices and K is the number of discovered keypoints. We use a learning rate of 0.001. The scalar loss coefficients (Section 3.2)  $\alpha_{kpt}$  and  $\alpha_{inf}$  are set to 1.0 and  $10^{-6}$  respectively. Figure 4 shows detailed description of network architectures used for the keypoint predictor  $\Phi$  and the influence predictor  $\Gamma$ .

**Datasets.** KeypointNet [7] dataset contains semantic 3D keypoint annotations for ShapeNet dataset [1]. Some models in KeypointNet are missing full keypoint annotations, therefore we use a subset of annotated keypoints that are contained in at least 80% of the models. KeypointNet also does not follow the standard training and testing splits from ShapeNet. We resample the KeypointNet dataset splits to make it compatible with the original ShapeNet splits.



Figure 2: Varying number of keypoints. The figure (a) shows the effect of different number of discovered keypoints (4, 8, 12, 16 from the top). The results are shown on randomly sampled results for ShapeNet Airplane category. Results in the table (b) are in terms of PCK@0.05 on airplanes from the KeypointNet dataset.

airplane																											
kŗ	ol kp2	kp3	kp4	kp5	kp6	kp7	kp8	best avg	kpl	kp2	kp3	kp4	kp5	kp6	6 kp7	kp8	best avg	kŗ	ol k	p2 k	p3 1	kp4	kp5	kp6	kp7	kp8	best avg
body 0.3	30 0.2	0.35	0.17	0.01	0.01	0.57	0.73		0.02	0.89	0.00	0.89	0.77	0.08	8 0.00	0.09		0.0	01 0.	02 0.	08 0	0.03	0.92	0.87	0.68	0.52	
wing 0.0	09 0.49	0.01	0.90	0.98	0.98	0.02	0.49		0.00	0.10	0.62	0.00	0.00	0.89	0.62	0.88		0.7	77 0.	96 0.	76 C	).95	0.00	0.01	0.02	0.03	
tail 0.1	10 0.0	0.00	0.01	0.00	0.00	0.82	0.01		0.65	5 0.00	0.00	0.02	0.00	0.01	0.00	0.01		0.0	00 0.	00 0.	00 0	0.01	0.00	0.00	0.12	0.86	
best 0.3	30 0.4	0.35	0.90	0.98	0.98	0.82	0.73	0.69	0.6	5 0.89	0.62	0.89	0.77	7 0.89	0.62	0.88	0.78	0.3	77 0.	<b>96</b> 0.	76 l	).95	0.92	0.87	0.68	0.86	0.85
		C	hen <i>et</i>	al.								Fern	andez	et al									ours				
								best					car				best										best
kj	pl kp	2 kp3	kp4	kp5	kp6	kp7	kp8	avg	kp	l kp2	kp3	kp4	kp5	5 kp(	5 kp7	kp8	avg	kj	pl k	:p2 k	р3	kp4	kp5	kp6	kp7	kp8	avg
roof 0.	00 0.0	0 0.00	0.00	0.00	0.00	0.00	0.00		0.0	0 0.50	0.01	0.00	) 0.00	0.0	0 0.00	0.00	)	0.	00 0.	.64 0.	.00 (	0.00	0.00	0.00	0.00	0.01	
wheels 0.	16 0.3	7 0.01	0.01	0.07	0.00	0.04	0.00		0.0	0 0.00	0.00	0.01	0.7:	5 0.0	1 0.02	0.58	3	0.	85 O.	.00 0.	.00	0.72	0.04	0.75	0.01	0.00	
body 0.	47 0.4	9 0.18	0.39	0.45	0.37	0.37	0.39	0.20	0.7	3 0.08	0.62	0.64	0.0	1 0.7	4 0.68	0.00		0.1	30 0.	.19 0.	.73 (	0.49	0.63	0.25	0.77	0.75	0.72
Dest 0.	4/ 0.4	9 0.18 C	0.39 'hen <i>et</i>	0.45 al	0.37	0.37	0.39	0.39	0.7	5 0.50	0.02	Eern	andez	$\frac{0.7}{2}$	4 0.00	0.50	0.00	0.0	85 0.	.04 0.	/3 0	0.72	0.03	0.75	0.77	0.75	0.75
		C	nen ei	uı.								1 CII	landez										ours	•			
													cha	ir													_4
	kp1	kp2	kp3	kp4	kp5	kp6	kp7	kp8	kp9 1	kp10 k	p11 kp	512	avg	kpl	kp2	kp3	kp4	kp5	kp6	kp7	kp8	kp9	kpl	0 kp	1 kp1	$\frac{2}{av}$	si g
bac	k 0.00	0.02	0.14	0.76	0.84	0.10	0.00	0.60	0.03	0.82 (	0.00 0.	.00		0.01	0.68	0.80	0.00	0.01	0.69	0.04	0.79	0.78	8 0.0	2 0.0	0 0.0	1	
sea	at 0.08	0.07	0.78	0.08	0.00	0.82	0.09	0.44	0.08	0.00	0.79 0.	.82		0.00	0.00	0.00	0.00	0.00	0.00	0.97	0.00	0.00	0.9	5 0.0	0 0.9	6	
leg	gs 0.77	0.80	0.09	0.02	0.00	0.05	0.76	0.04	0.79	0.00 (	0.09 0.	.08		0.76	0.00	0.00	0.75	0.75	0.00	0.01	0.01	0.00	) 0.0	1 0.7	6 0.0	1	r i
bes	st 0.77	0.80	0.78	0.76	0.84	0.82	0.76	0.60	0.79	0.82 (	0.79 0.	.82 (	0.78	0.76	0.68	0.80	0.75	0.75	0.69	0.97	0.79	0.70	8 0.9	5 0.7	6 0.9	6 0.8	80
						Cne	en et a	l.											ге	rnand hest	ez ei	aı.					
								kp l	kp2	kp3 1	kp4 k	p5 1	kp6	kp7	kp8 l	cp9 k	.p10 k	pll k	cp12	avg							
							back	0.00	0.74	0.00	0.73 0	.00 (	).01 (	0.01	0.88 (	0.01	).94 0	0.96 (	0.01								
							seat	0.86	0.61	0.02	).67 <b>0</b>	.01 (	).88 (	0.01	0.01 (	0.02 (	0.00 0	0.00	0.91								
							legs	0.01	0.30	0.94 (	0.31  0.000	.93 (	).52 (	0.88	0.00	0.91	0.01 0	0.01	0.48	0.00							
							best	0.86	0.74	0.94 (	). 73 0.	.93 (	01 88 0	0.88 Irs	0.88 [ (	0.91 (	0.94 0	0.96 (	9.91	0.88							
													00	• ,		<i>.</i> .											
									(a) 1	insur	ervis	еа к	eypo	ints (	correl	atior	1										
			airp	lane	cap	ca	r ch	air g	guitar	knif	e lap	top	mote	orbike	mu	g sk	atebo	ard	table	pist	ol	bag	rocl	cet e	arpho	one	lamp
Chen et	t al. [2	]	0.	69	0.24	0.3	9 0.	78	0.97	0.94	0.9	95	0.	91	0.5	)	0.89		0.75	0.7	8 (	0.35	0.5	6	0.30	)	0.50
Fernance ours	dez et	al. [3]	0.	78 85	0.45 <b>0.71</b>	0.6 <b>0.7</b>	60. 30.	80 88	0.93 <b>0.99</b>	0.92 <b>0.96</b>	2 0.8 5 <b>0.</b> 9	85 96	0. <b>0.</b>	90 <b>93</b>	0.73 <b>0.9</b> 4	3 1	0.92 <b>0.96</b>		0.85 0.92	0.6 <b>0.9</b>	0 ( 1 (	).72 ).85	0.6 <b>0.9</b>	01 0	0.24 0.72	+ 2	0.40 <b>0.53</b>

(b) average unsupervised keypoints correlation

Figure 3: Semantic part correspondence. We evaluate semantic part correspondence for the ShapeNet Car category. The tables (a) shows the frequency of each unsupervised keypoint  $[kp^*]$  being associated with a given object part. Chen *et al.* [2] show worse performance in this task because that methods tends to predict keypoints inside the object far from the annotated object surface. We also report the average unsupervised keypoints correlation for each category (b).  $\uparrow$  is better.

Operation	Kernel	Output channels	Output size	Activation
nput x	-	3	1024	-
Conv 1D	1	64	1024	ReLU
Conv 1D	1	128	1024	ReLU
Conv 1D	1	256	1024	None
Max. pool	1024	256	1	-
Squeeze	-	256	-	-
Linear	-	256	-	LReLU
Linear	-	512	-	LReLU
Linear	-	256	-	LReLU
Linear	-	3 <i>·K</i>	-	None
Reshape	-	3	K	-

(a) Keypoint predictor  $\Phi$ 

(b) Influence predictor  $\Gamma$ 

Figure 4: Network architectures. The network architectures are based on a PointNet encoder [5, 6]. K is the number of discovered keypoints, C is the number of cage vertices. LReLU stands for Leaky ReLU with 0.1 negative slope.



Figure 5: Shape control via unsupervised 3D keypoints. We show steps from shape editing using our unsupervised 3D keypoints. Please refer to our supplementary video on our project page to see the editing in action.



Figure 6: Unsupervised 3D keypoints on real-world 3D scans. Randomly sampled results with 8 unsupervised keypoints on real-world 3D scans of shoes from Google Scanned Objects dataset [4].



Figure 7: Unsupervised 3D keypoints. Randomly sampled results with 8 unsupervised keypoints for ShapeNet Airplane category.



Figure 8: Unsupervised 3D keypoints. Randomly sampled results with 8 unsupervised keypoints for ShapeNet Guitar category.



Figure 9: Unsupervised 3D keypoints. Randomly sampled results with 12 unsupervised keypoints for ShapeNet Chair category.



Figure 10: Unsupervised 3D keypoints. Randomly sampled results with 8 unsupervised keypoints for ShapeNet Car category.

## References

- 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. 2
- [2] Nenglun Chen, Lingjie Liu, Zhiming Cui, Runnan Chen, Duygu Ceylan, Changhe Tu, and Wenping Wang. Unsupervised learning of intrinsic structural representation points. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9121–9130, 2020. 4
- [3] Clara Fernandez-Labrador, Ajad Chhatkuli, Danda Pani Paudel, Jose J Guerrero, Cédric Demonceaux, and Luc Van Gool. Unsupervised learning of category-specific symmetric 3d keypoints from point sets. *European Conference on Computer Vision (ECCV)*, 2020. 4
- [4] GoogleResearch. Google scanned objects. https://fuel.ignitionrobotics.org/1.0/GoogleResearch/fuel/ collections/Google%20Scanned%20Objects, September 2020. 6
- [5] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In Proceedings of the IEEE conference on computer vision and pattern recognition, pages 652–660, 2017. 5
- [6] Wang Yifan, Noam Aigerman, Vladimir G Kim, Siddhartha Chaudhuri, and Olga Sorkine-Hornung. Neural cages for detail-preserving 3d deformations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 75–83, 2020. 5
- [7] Yang You, Yujing Lou, Chengkun Li, Zhoujun Cheng, Liangwei Li, Lizhuang Ma, Cewu Lu, and Weiming Wang. Keypointnet: A large-scale 3d keypoint dataset aggregated from numerous human annotations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13647–13656, 2020. 2