# Skeleton Merger: an Unsupervised Aligned Keypoint Detector

Ruoxi Shi, Zhengrong Xue, Yang You, Cewu Lu*

Shanghai Jiao Tong University

`eliphat, xuezhengrong, qq456cvb, lucewu@sjtu.edu.cn`

## Abstract

*Detecting aligned 3D keypoints is essential under many scenarios such as object tracking, shape retrieval and robotics. However, it is generally hard to prepare a high-quality dataset for all types of objects due to the ambiguity of keypoint itself. Meanwhile, current unsupervised detectors are unable to generate aligned keypoints with good coverage. In this paper, we propose an unsupervised aligned keypoint detector, Skeleton Merger, which utilizes skeletons to reconstruct objects. It is based on an Autoencoder architecture. The encoder proposes keypoints and predicts activation strengths of edges between keypoints. The decoder performs uniform sampling on the skeleton and refines it into small point clouds with pointwise offsets. Then the activation strengths are applied and the sub-clouds are merged. Composite Chamfer Distance (CCD) is proposed as a distance between the input point cloud and the reconstruction composed of sub-clouds masked by activation strengths. We demonstrate that Skeleton Merger is capable of detecting semantically-rich salient keypoints with good alignment, and shows comparable performance to supervised methods on the KeypointNet dataset. It is also shown that the detector is robust to noise and subsampling. Our code is available at https://github.com/eliphatfs/SkeletonMerger.*

## 1. Introduction

Being able to fully understand an object is arguably the ultimate goal of computer vision. For 3D point clouds, detecting semantic keypoints is currently the most promising and widely adopted approach. [25, 21] Keypoints are crucial to the success of many vision applications such as object tracking, shape registration and in robotics [15, 24, 3, 5]. In many actual cases where objects from the same category are compared, we desire keypoints to be not only accurately located but also aligned within a certain category for performing high-level vision tasks such as 3D object recog-
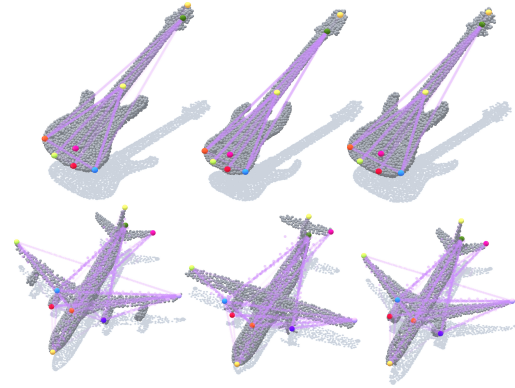


Figure 1. **Examples of detected aligned keypoints and the corresponding skeletons from our detector.** Skeleton is shown in purple lines. Objects are from ModelNet40 [27] and ShapeNet-CoreV2 [4] datasets.

nition and reconstruction. [29]

However, to decide whether a point contains semantics requires high-level intelligence, as 'semantics' itself is ambiguously defined. Different people would present different understandings of semantic points. Therefore, very limited human annotated quality data are available so far [31]. Consequently, supervised methods can only deal with very limited range of objects covered in the datasets despite their success on many other tasks.

Most unsupervised methods, either traditional hand-crafted ones [19, 33] or deep learning-based [12], take advantage of geometric properties to detect keypoints. While stable, these keypoints are often not rich in semantics, and the coverage of keypoints on the input point cloud is generally low, especially under a small number of points, which limits their performance in downstream tasks. Moreover, the keypoints detected are neither ordered nor aligned. A very recent approach [8] can learn aligned 3D keypoints by decomposing keypoint coordinates into a low-rank non-rigid shape representation. However, in categories like the airplane where objects do not necessarily share highly similar geometric shapes, its performance dramatically declines.

In view of the challenges above, we propose Skeleton Merger, an unsupervised keypoint detector that can extract

---

*Cewu Lu is corresponding author. He is also the member of Qing Yuan Research Institute and MoE Key Lab of Artificial Intelligence, AI Institute, Shanghai Jiao Tong University, China

Figure 2. **Describe objects with skeletons.** (a) We observe that human vision can easily distinguish objects with skeletons, but it is harder with keypoints only. (b) Histogram of nearest neighbour distances of points in point clouds to its skeleton (red), large quantity of (3200) uniformly sampled points (green) in the bounding box and keypoints (blue).

salient and aligned keypoints. As its name suggests, Skeleton Merger attempts to reconstruct a point cloud from its skeleton. Some examples of detected keypoints and the skeletons are shown in Fig. 1.

Utilizing skeleton to represent a point cloud is inspired by the skeleton extraction problem [6] in traditional graphics. We hold the belief that skeleton is a better presentation method than that adopted by current deep learning-based frameworks both intuitively and statistically. Qualitatively speaking, we observe that human vision tends to use a 'joint-skeleton' cognitive pattern to recognize things, as is shown in Fig. 2 (a). While accurate and descriptive semantics is mostly provided by keypoints, with some discrete keypoints as joints alone, it is merely impossible for human beings to recognize objects. However, if we introduce some auxiliary line segments connecting certain keypoints together to form a skeleton, humans can easily distinguish the object. Quantitatively speaking, statistical results in Fig. 2 (b) show that points of a point cloud are in general closer to those of a skeleton (red), compared with keypoints (blue) or uniform sampling inside the bounding box (green). This indicates skeleton can better fit with local shape features of the original point cloud.

On implementation level, our Skeleton Merger follows a deep Autoencoder architecture. The encoder generates keypoint proposals. It also predicts activation strengths of edges between keypoints. The decoder generates a skeleton from these results, which is essentially a graph of keypoints. It performs uniform sampling on each edge, adds distinct activation strength and offset to refine the shape. In this way, each skeleton edge is essentially a small subcloud. The decoder finally merges them altogether to form a reconstruction point cloud. Noticing that the order of skeleton edges is predefined by the encoder, the alignment of keypoints is therefore considerably improved.

The crucial problem now becomes how to construct a loss function that can evaluate how well the refined skeleton reconstructs the original point cloud. Following the idea of traditional Chamfer loss [2], which is the sum of forward and backward losses, we come up with Composite Chamfer Distance (CCD), which is the sum of fidelity (forward) and coverage (backward) losses. CCD measures the distance between the input and the reconstructed point cloud composed of many sub-clouds masked by activation strengths.

Experimental results show that Skeleton Merger is capable of detecting semantically-rich salient keypoints with good alignment. Our detector achieves significantly better performance than current unsupervised detectors on the KeypointNet [31] dataset. In fact, its performance is even comparable to supervised methods. Results also show that our detector is robust to noise and subsampling.

## 2. Related work

**Curve skeleton extraction** In the context of computer graphics, skeleton refers to the 1D structure which is a simplified representation of the geometry of a 3D object. Cornea *et al*. [6] provided an overview of curve skeleton algorithms [1, 14, 18] and their applications. In our paper, however, skeleton refers to a graph of keypoints that represents topology of the object. The purpose of our skeleton is not only to provide a rough geometric shape of the original object, but to help improve alignment of keypoints.

**Deep learning on 3D point clouds** Currently, various deep learning-based techniques that consume point clouds [13, 16, 17, 28, 26, 11] have been developed. They initially aim at basic tasks such as classification and segmentation, but can be adapted for more high-level 3D perception tasks like point cloud registration, pose estimation and 3D reconstruction [9, 7, 32, 10]. PointNet [16] proposed by Qi *et al*. is a pioneering work that first enables neural networks to directly process raw point cloud data. In PointNet, the input points pass through per-point multi-layer perceptron and a symmetric max pooling operation, ending up with a global feature. The global feature is then used for various downstream tasks. However, PointNet only takes notice of the global information, neglecting local details. Thus, Qi *et al*. extended PointNet to PointNet++ [17], where PointNet is applied hierarchically on different spatial scales for better performance on point clouds. Our encoder utilizes PointNet++ as a point cloud processing block.

**Unsupervised 3D keypoint detectors** Currently, most 3D keypoint detectors remain to be hand-crafted. Popular hand-crafted detectors such as Harris 3D [19], ISS [33], HKS [20], SHOT [23] take advantage of geometric properties to select most salient keypoints. As geometric characteristics are quite complex in 3D objects, most keypoints detected are neither semantically salient nor well aligned.

To the best of our knowledge, USIP [12] is the first learning-based 3D keypoint detector. USIP takes advantage
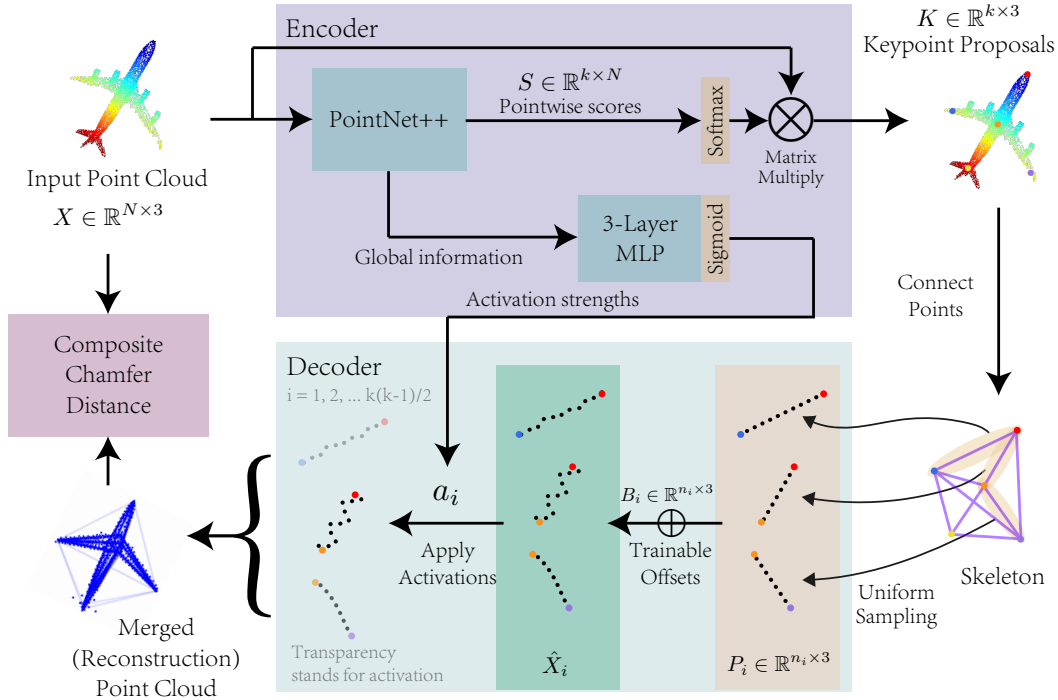
Figure 3. **Pipeline of Skeleton Merger.** The encoder takes $n$ points as input and utilizes a PointNet++ network to generate keypoint proposals and activation strengths. 'MLP' stands for multi-layer perceptron. Batch-norm and ReLU are used for the MLP. The decoder samples on edges of the skeleton, refines them by adding positional offsets that are directly optimized as parameters, and merges the refined skeleton edges with the activation strengths. Composite Chamfer Distance is applied between the reconstruction result and the input point cloud.

of probabilistic chamfer loss that greatly enhances repeatability. However, while keypoints detected are stable under geometric transformations, they may have poor coverage on the point cloud, especially when the number of keypoints is small. Furthermore, keypoints detected are unordered and not aligned. Very recently, Fernandez *et al*. [8] propose an unsupervised approach that can learn aligned 3D keypoints through decomposing keypoint coordinates into a low-rank non-rigid shape representation. This leads to degraded performance on classes where objects are not necessarily similar in geometry such as guitars and airplanes of irregular shapes. Moreover, all shapes are assumed to be axisymmetric in their method, while our detector needs no symmetry prior, so it can be applied to a wider range of objects.

## 3. Methods

Now we propose the unsupervised aligned keypoint detector, Skeleton Merger, which is based on a deep Autoencoder framework. The full pipeline is shown in Fig. 3.

The encoder first proposes an ordered list of $k$ keypoints $K \in \mathbb{R}^{k \times 3}$. Connecting each pair of keypoints, a rough skeleton of the input point cloud is generated, composed of $C_k^2 = k(k-1)/2$ edges. The skeleton is further refined into $k(k-1)/2$ point clouds through uniform sampling and addi-

tion with a sequence of learned offset vectors. These clouds are masked by $k(k-1)/2$ activation strengths, which also come from the encoder, and merged into a final reconstruction point cloud at the output of the decoder. The Composite Chamfer Distance (CCD) is applied between the merged reconstruction point cloud and input point cloud to guide the unsupervised learning process.

In the following subsections, we will introduce in detail each module in the pipeline.

### 3.1. Skeleton and activation strengths

The key of the Skeleton Merger architecture lies in the utilization of the skeleton to reconstruct objects.

Directly reconstructing the original point cloud from keypoints is essentially transforming a point set into a larger one, which by its nature is unordered. Furthermore, keypoints only contain sparse semantic information, and experiments show that reconstruction by this method may lead to poor coverage of the original point cloud. In contrast, the skeleton provides a basic geometric shape of the original object, instead of sparse semantics. Meantime, it can be easily constructed with alignment by connecting pairs of keypoints in order.

There are edges in the skeleton, however, that do not ac-

tually exist in the original point cloud. They will bring unwanted noise to the reconstruction. A mechanism is thus required for masking out these non-existing edges. So together with the skeleton, a list of activation strength values are predicted by the model. These values indicate existence of edges in the skeleton. This step requires the model to learn the skeleton topology of the reconstruction target object, which, combined with ordered construction of the skeleton edges, improves alignment of keypoints.

### 3.2. Encoder, the keypoint proposer

In the encoder, a PointNet++ [17] network is first applied to obtain $k$ pointwise scores $S \in \mathbb{R}^{k \times N}$ of the input point cloud. The pointwise scores are then activated by a softmax function. A weighted average of the input point cloud $X \in \mathbb{R}^{N \times 3}$ is computed from $\texttt{softmax}(S)$ and $X$ to form the final keypoint proposals $K \in \mathbb{R}^{k \times 3}$. The weighted average is implemented by a matrix multiplication, shown in Fig. 3.

Besides $k$ keypoints, the existence of $k(k-1)/2$ edges between each pair of keypoints is encoded as well. A 3-layer MLP (multi-layer perceptron) accepts the *global feature vector* generated by the PointNet++ network and predicts $k(k-1)/2$ sigmoid-activated activation strengths $a \in \mathbb{R}^{k(k-1)/2}$. The activation strengths are used to mask skeleton edges before the merging stage in decoder.

### 3.3. Decoder, the skeleton refiner

The decoder takes as input the keypoints $K \in \mathbb{R}^{k \times 3}$ proposed by the encoder. A uniform sampling operation is performed on the edges to get $k(k-1)/2$ small point clouds, $P_1, P_2, \ldots, P_{k(k-1)/2}$. The number of points sampled $n_i$ is in proportion to the length of each edge.

For each small point cloud $P_i \in \mathbb{R}^{n_i \times 3}$, pointwise position offset $B_i \in \mathbb{R}^{n_i \times 3}$ is added to the initial points sampled as a refinement to the skeleton formed by straight lines. $B_i$'s are directly optimized as parameters of the network. In order to keep the refinement localized, a Ridge (L2) regularization is imposed on the learnt position offsets.

The $k(k-1)/2$ refined point clouds $\hat{X}_i = P_i + B_i$ are merged into a single point cloud with the activation strengths $a_i$ from the encoder. The CCD is then applied between the reconstruction point cloud and the input point cloud $X$ as the loss to guide the training.

### 3.4. Composite Chamfer Distance

To put the reconstruction process into practice, it is essential to establish a loss function between the input point cloud $X$ and the reconstruction point cloud composed of different parts $\hat{X}_i$ masked by activation strengths $a_i$. We generalize the Chamfer Distance to take into account the activation strengths, and proposes the Composite Chamfer Distance (CCD).

Similar to the regular Chamfer Distance, CCD is a sum of fidelity (forward) and coverage (backward) losses. However, the reconstruction result is composed of several sub-clouds, while the input is simply one large point cloud. This asymmetry leads to asymmetry in designs of fidelity and coverage losses, making it to fit the nature of the problem.

The fidelity loss is a straightforward extension to the Chamfer Distance where the activation strength is applied to each sub-cloud, as shown in Eq. 1:

$$\mathcal{L}_f = \sum_i a_i \sum_{\hat{p} \in \hat{X}_i} \min_{p_0 \in X} ||\hat{p} - p_0||_2 . \qquad (1)$$

Then it comes to the coverage loss. If we do the same simple extension to the Chamfer Distance, activation strengths $a_i$ will go to zero when the loss gets minimized, which prevents the model from learning anything meaningful. The problem is that, a reconstruction point with a small $a_i$ value does not contribute to coverage as much as one with a large $a_i$ value. Therefore, more points than only the one with minimal distance should be considered if its activation strength is not large enough.

In view of this, we come up with the following coverage loss, which involves point-wise sorting of sub-clouds and weighted averaging instead of a simple minimum. The algorithm to generate the coverage loss is shown in Alg. 1. An illustrating example is shown in Fig. 4.

---

**Algorithm 1** Coverage loss of CCD

---

    **Input**: $X, \hat{X}_1 \ldots \hat{X}_{k(k-1)/2}, a_1 \ldots a_{k(k-1)/2}$
    **Parameter**: $\gamma$
    **Output**: $\mathcal{L}_c$
1: **for** $p_0 \in X$ **do**
2:     $\mathcal{R} \leftarrow \bigcup \{\hat{X}_1 \ldots \hat{X}_{k(k-1)/2}\}$
3:     $w \leftarrow 0$
4:     **while** $w < 1$ and $\mathcal{R} \neq \emptyset$ **do**
5:         Find $\hat{p}$ by $\min_{\hat{p} \in \mathcal{R}} ||\hat{p} - p_0||_2$
6:         Find $i$ with $\hat{p} \in \hat{X}_i$
7:         $\mathcal{L}_c \leftarrow \mathcal{L}_c + a_i ||\hat{p} - p_0||_2$
8:         $w \leftarrow w + a_i, \mathcal{R} \leftarrow \mathcal{R} \setminus \hat{X}_i$
9:     **end while**
10:    **if** $w < 1$ **then**
11:        $\mathcal{L}_c \leftarrow \mathcal{L}_c + \gamma(1 - w)$
12:    **end if**
13: **end for**

---

Each point in the input point cloud is treated separately (the outer 'for' loop starting at line 1) with an iterative process: the point with minimal distance in the collection of sub-clouds is selected (line 5 and 6), and the distance multiplied with the activation strength is the contribution of the current iteration to the coverage loss (line 7), after which the entire sub-cloud is removed from the collection under
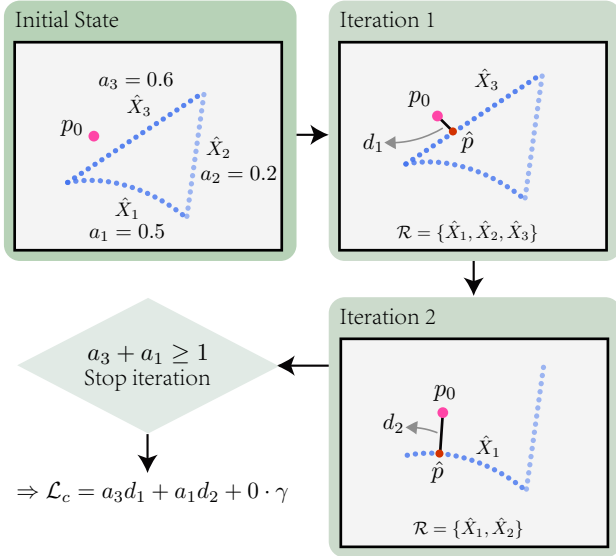
For each point $p_0$ in the input point cloud $X$



Figure 4. **An illustrating example of the coverage loss algorithm.** The input cloud is a curved poly-line in the example. Currently the focus $p_0$ is the marked pink point. In each iteration, the nearest point of $p_0$ in $\mathcal{R}$ is selected and the distance is computed. Then the small point cloud that contains the selected point is discarded from $\mathcal{R}$ and the next iteration starts. The iteration stops when selected $a_i$'s sum up to 1 (as shown), or when $\mathcal{R}$ is empty, in which case $\gamma$ would be applied as a penalty for the rest of $w$. See the text for more details about coverage loss.

consideration (line 8). The iteration stops after activation strengths sum up to 1 (notice the $w < 1$ condition in the while loop starting at line 4). If it never reaches 1, a high penalty of $\gamma$ is imposed (line 10 to 12). $\gamma$ is set to 20.0 in the experiments.

Intuitively, sub-clouds that contribute more to the coverage have small backward Chamfer Distance, and due to the cutting mechanism, the activation strengths of these parts will get larger. The fidelity loss, in contrast, reduces the activation strengths of non-existing edges in the skeleton. By applying CCD as loss, the network is forced to generate a set of sub-clouds with reasonable activation values, which enables the training process.

Meantime, alignment is assured since the activation-generating MLP only sees global information about the point cloud. The MLP learns essentially the topology of target model skeletons, and a wrong ordering of keypoints will lead to high fidelity and coverage losses.

The final CCD loss is a weighted sum of the two parts, fidelity loss $\mathcal{L}_f$ and coverage loss $\mathcal{L}_c$, as shown in Eq. 2. The weight coefficients $\lambda_f$ and $\lambda_c$ can be tuned. They default to be the same in the experiments.

$$\mathcal{L} = \lambda_f \mathcal{L}_f + \lambda_c \mathcal{L}_c. \qquad (2)$$

## 4. Results

In this section, we evaluate keypoint saliency, alignment and repeatability of the proposed Skeleton Merger, and give some qualitative results of the keypoint detector.

### 4.1. Metric for alignment with human annotations

Currently, there lacks a metric for evaluating correspondence between two sets of semantic labels for keypoints, where a consistent label is given to keypoints with the same semantics in each set, but the relation between labels of two sets is unknown. We propose Dual Alignment Score (DAS) metric for this evaluation.

**Dual Alignment Score** In order to evaluate whether our keypoints are consistent in semantics, we propose Dual Alignment Score (DAS). In DAS calculation, a reference point cloud is used for semantic assignment and another model is used for the actual evaluation. On the reference model, we assign each predicted keypoint with the semantic index of the closest point from human annotations. Then on the other point cloud, the corresponding predicted point is found since our keypoints are aligned. A score is calculated by the mean accuracy whether the closest human annotated keypoint of this point is of the same semantic index. In the opposite direction, order of our predicted keypoints is used to assign semantic labels for human annotations, and the process is repeated. By averaging scores in these two directions we get the Dual Alignment Score (DAS). Fig. 5 shows an illustration of the DAS computation.
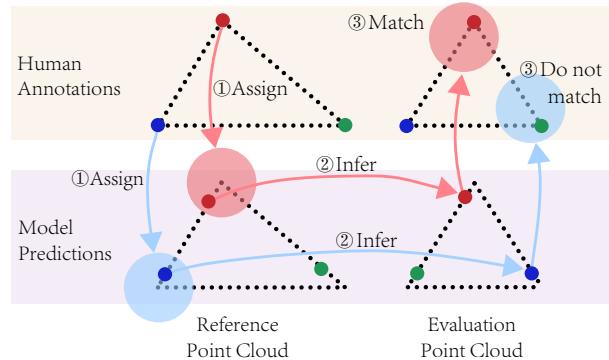


Figure 5. **An illustration of DAS computation process.** Only one matching direction from predicted semantic labels to annotations is shown. We first assign corresponding labels by finding the nearest neighbour of predicted keypoints on annotations of the reference point cloud. Then the human annotated keypoints with same semantic labels are inferred for the evaluation point cloud. Finally, we see if the nearest neighbour of the inferred annotations are of same predicted labels, and compute an accuracy score.
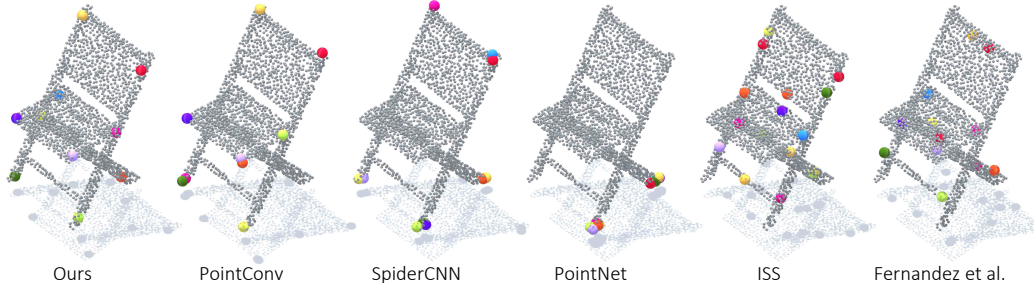
5

Figure 6. **Comparison between different supervised and unsupervised detectors.** Our detector produces salient, semantically rich keypoints and is comparable to or even outperforms supervised methods, as shown.

## 4.2. Evaluation on human annotated data

We first evaluate saliency and alignment of keypoints detected by the proposed Skeleton Merger. More specifically, we consider evaluation on two large-scale keypoint datasets, KeypointNet[31] and data from SyncSpecCNN [30], where keypoints are annotated by experts with semantic correspondence labels.

**Keypoint saliency**  For keypoint saliency, we compute the mean Intersection over Unions (mIoU) [22] metric, where

$$IoU = \frac{TP}{TP + FP + FN}. \qquad (3)$$

mIoU is computed with a threshold of 0.1 in euclidean distance. Skeleton Merger is trained and tested on KeypointNet. Also, several supervised networks, PointNet [16], SpiderCNN [28] and PointConv [26], are trained on KeypointNet to predict the likeliness of each point to be a keypoint as a comparison. In addition, the recent unsupervised detector [8] by Fernandez *et al.*[1] and two traditional methods, Harris3D [19] and ISS [33], are also compared.

The results are shown in Tab. 1. It can be seen that in these categories where keypoints are (at least partially) well defined, our unsupervised detector shows competitive performance to, or even outperforms, the supervised networks, in terms of mIoU, and is by far superior to traditional methods in detecting salient keypoints. A visualized comparison between different methods is shown in Fig. 6.

**Keypoint alignment**  For keypoint alignment, Skeleton Merger is trained on the ShapeNetCoreV2 [4] dataset, and DAS described in the previous section is evaluated on the KeypointNet[31] and SyncSpecCNN [30] datasets. The results are shown in Tab. 2. We compare the DAS scores with the method of Fernandez *et al.* [8].

---

[1]The method in [8] requires a category-specific symmetry prior, which is not available in the Guitar class.

|                   | Airplanes | Chairs | Guitars |
|-------------------|-----------|--------|---------|
| PointNet          | 45.4      | 23.8   | 0.2     |
| SpiderCNN         | 55.0      | 49.0   | 17.0    |
| PointConv         | 93.5      | 86.0   | 84.9    |
| Harris3D          | 42.8      | 15.1   | 33.1    |
| ISS               | 36.3      | 11.6   | 37.0    |
| Fernandez *et al*.| 69.7      | 51.2   | -       |
| Ours              | 79.4      | 68.4   | 55.0    |

Table 1. mIoU scores of Skeleton Merger, different supervised networks, method of Fernandez *et al*. and traditional keypoint detectors on KeypointNet.

|                    | Airplane (K) | Chair (K) | Chair (S) |
|--------------------|--------------|-----------|-----------|
| Fernandez *et al*. | 61.4         | 64.3      | 54.2      |
| Ours               | **77.7**     | **76.8**  | **73.8**  |

Table 2. DAS scores of Skeleton Merger and method of Fernandez *et al*. on KeypointNet (K) and SyncSpecCNN (S).

## 4.3. Repeatability

In this section we investigate the repeatability of detected keypoints to Gaussian additive noise and point cloud downsampling. Gaussian noises of different strengths and different sampling ratios are applied on the point clouds, and the same Skeleton Merger network trained on ShapeNetCoreV2 [4] is applied for keypoint detection.

The keypoints that are detected on these modified point clouds are compared with those detected from the clean, original point cloud *in order*, that is, keypoints are compared one-by-one as a list instead of a set as in most previous works. If the distance between a keypoint detected from the original point cloud and a keypoint detected under noise or subsampling is smaller than 10% of the model size, the keypoint is considered repeatable.

The results are shown in Fig. 8. We compare the results with ISS [33]. It can be seen that the aligned keypoints

|  | mIoU | | DAS | |
| --- | --- | --- | --- | --- |
|  | Airplane | Chair | Airplane | Chair |
| Full Skeleton Merger | 79.4 | **68.4** | **77.7** | **76.8** |
| No activation strengths | 55.5 | 8.4 | 72.1 | 65.2 |
| No fidelity loss | 78.2 | 60.0 | 76.2 | 74.4 |
| No coverage loss | 17.8 | 1.1 | 35.6 | 37.3 |
| No offsets | **85.6** | 62.0 | 72.4 | 75.6 |

Table 3. Ablation study of different components in Skeleton Merger.



Full Skeleton Merger     No activation strengths     No fidelity loss     No coverage loss     No offsets

Figure 7. **Visualization results of the ablation study.** Minor degeneracies can be seen in models without offsets or the fidelity loss. Major performance drop is seen if either the activation strengths or the coverage loss is removed.

detected remain highly repeatable under Gaussian noise or downsampling. They also stay well-aligned.
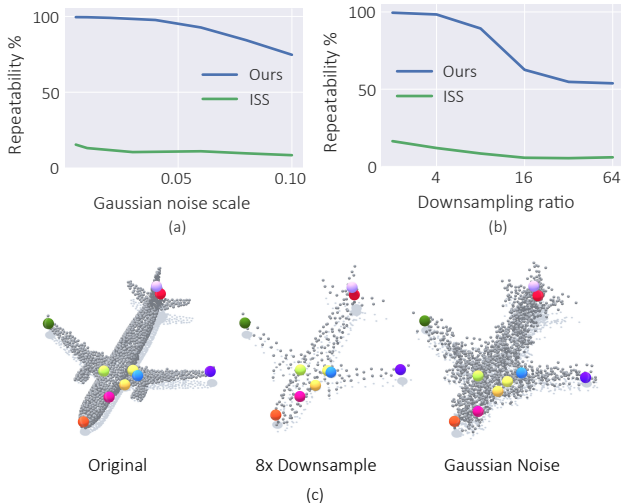


Figure 8. **Repeatability of Skeleton Merger and ISS on ShapeNetCoreV2.** (a) Repeatability under Gaussian noise. (b) Repeatability under downsampling. (c) Visualization results of keypoints from Skeleton Merger under the scenarios. Downsample rate is 8x and the Gaussian noise scale is 0.05.

## 4.4. Ablation study

We carried out experiments on the effectiveness of proposed components in Skeleton Merger. Tab. 3 and Fig. 7 shows the results from the network with different settings. The models are trained on ShapeNetCoreV2 [4] and evaluated on KeypointNet [31].

**Activation strengths** Without activation strengths, *i.e.* all activations of sub-clouds are set to 1, the network has to cope with the non-existing edges to fit the input point clouds, and there lacks a mechanism to enforce alignment of keypoints (the skeleton topology is no longer utilized for alignment), thus the full model outperforms model of this version both in saliency (mIoU) and alignment (DAS).

**Trainable offsets** Without offsets, only straight lines between keypoints are allowed. These lines fit the shape of airplanes well, and removing offsets reduces variance in the model, so this version outperforms the full model in terms of keypoint saliency. In other cases, however, the model without offsets is not so lucky. It suffers from filling squares, balls and other shapes with straight lines.

**Composite Chamfer Distance** The Composite Chamfer Distance is at the core of the training process.

Without the fidelity loss, activation strengths soon go to 1 because of the stop-iteration procedure in the coverage loss. As a result, the network can only learn meaningful latents in the first several epochs, and stops improving due to the same reasons without activation strengths, causing degeneracies in performance.

Without the coverage loss, activation strengths soon go to 0 as fidelity loss is minimized, preventing the network from learning anything meaningful. This also emphasizes the importance of coverage in keypoint detection.
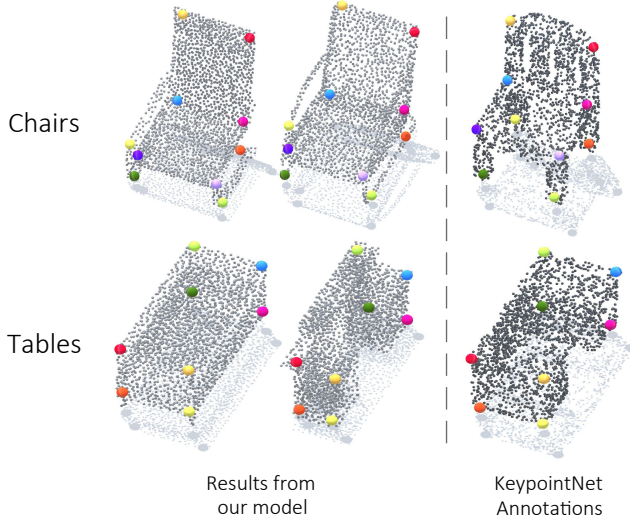
Figure 9. **Examples of keypoints detected by Skeleton Merger on ShapeNetCoreV2.** Keypoints detected on chairs and tables, together with a set of annotations from KeypointNet are displayed.
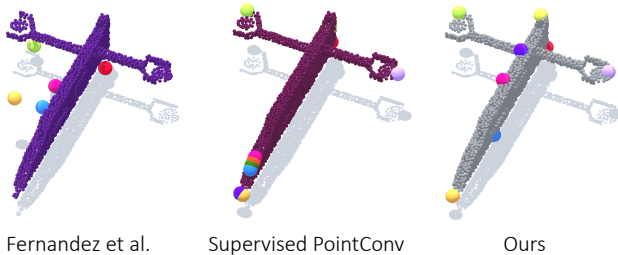


Figure 10. **Comparison of different detectors on an irregular-shaped airplane.** The method from Fernandez *et al.* fails to generalize on this instance, and supervised PointConv shows degeneracies, while our method works well.

## 4.5. Qualitative results and limitations

In this section, we give some qualitative results of our method and discuss the keypoint ambiguity and limitations.
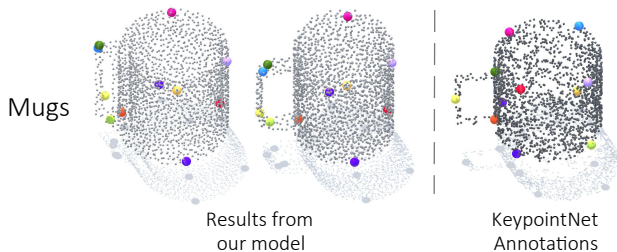


Figure 11. **Examples of keypoint detections and human annotations in the mug category.** The points on the lip of the mugs share the same semantics, which leads to ambiguity in keypoint definition.

**Qualitative visualization results** Fig. 9 shows some visualization results of keypoints detected by Skeleton Merger on different object categories. It can be seen that these points are well-aligned between different instances of objects and cover most points with semantics in the point cloud. They correspond well to the human annotations.

**Irregular-shaped instances** It is worth mentioning the generalization capability of our keypoint detector to irregular-shaped objects in an category. As shown in Fig. 10, the model of Fernandez *et al.* [8] fails to generalize to the irregular-shaped airplane, and the supervised PointConv [26] network shows some minor degeneracies in keypoint detection due to low frequencies of irregular objects appearing in the training set. Our method still works fine on this irregular-shaped instance.

**Keypoint ambiguity** It is demonstrated in Fig. 11 that keypoint definitions are ambiguous in some objects. Points on the lip of the mug, for example, are equivalent in semantics due to circular symmetry of the shape. Skeleton Merger and the KeypointNet ground truth both yield symmetric points, but points with different angles are selected.

The ambiguity makes it hard to aggregate human annotations to obtain a high-quality dataset for a wide range of objects, such as jars and cameras. As discussed before, this imposes a strong limit on the application of supervised methods for keypoint detection.

**Limitations** Skeleton Merger is capable of generating semantically-rich and well-aligned keypoints. However, it is less sensitive to local semantics than global coverage. For example, joints are already covered by a cross of two skeleton edges. Selecting keypoints at these points may not reduce the global losses.

## 5. Conclusion

In this paper, we present Skeleton Merger, an unsupervised aligned keypoint detector. Composite Chamfer Distance (CCD) is proposed as a loss function to guide the network to detect high-quality keypoints by reconstructing a point cloud through refining its skeleton. Evaluations are performed on the quality and repeatability of detected keypoints. Our detector shows impressive performance detecting salient and well-aligned keypoints.

## Acknowledgement

# References

[1] Oscar Kin-Chung Au, Chiew-Lan Tai, Hung-Kuo Chu, Daniel Cohen-Or, and Tong-Yee Lee. Skeleton extraction by mesh contraction. *ACM transactions on graphics (TOG)*, 27(3):1–10, 2008. 2

[2] Harry G Barrow, Jay M Tenenbaum, Robert C Bolles, and Helen Cf Wolf. Parametric correspondence and chamfer matching: Two new techniques for image matching. In *Proceedings: Image Understanding Workshop*, pages 21–27. Science Applications, Inc Arlington, VA, 1977. 2

[3] 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 586–606. International Society for Optics and Photonics, 1992. 1

[4] Angel X Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, Jianxiong Xiao, Li Yi, and Fisher Yu. ShapeNet: An Information-Rich 3D Model Repository. Technical Report arXiv:1512.03012 [cs.GR], Stanford University — Princeton University — Toyota Technological Institute at Chicago, 2015. 1, 6, 7

[5] Changhyun Choi and Henrik I Christensen. Real-time 3d model-based tracking using edge and keypoint features for robotic manipulation. In *2010 IEEE International Conference on Robotics and Automation*, pages 4048–4055. IEEE, 2010. 1

[6] Nicu D Cornea, Deborah Silver, and Patrick Min. Curve-skeleton properties, applications, and algorithms. *IEEE Transactions on visualization and computer graphics*, 13(3):530, 2007. 2

[7] Gil Elbaz, Tamar Avraham, and Anath Fischer. 3d point cloud registration for localization using a deep neural network auto-encoder. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4631–4640, 2017. 2

[8] 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. *arXiv preprint arXiv:2003.07619*, 2020. 1, 3, 6, 8

[9] Yulan Guo, Hanyun Wang, Qingyong Hu, Hao Liu, Li Liu, and Mohammed Bennamoun. Deep learning for 3d point clouds: A survey. *IEEE transactions on pattern analysis and machine intelligence*, 2020. 2

[10] Xianfeng Han, Hamid Laga, and Mohammed Bennamoun. Image-based 3d object reconstruction: State-of-the-art and trends in the deep learning era. *IEEE transactions on pattern analysis and machine intelligence*, 2019. 2

[11] Asako Kanezaki, Yasuyuki Matsushita, and Yoshifumi Nishida. Rotationnet: Joint object categorization and pose estimation using multiviews from unsupervised viewpoints. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5010–5019, 2018. 2

[12] Jiaxin Li and Gim Hee Lee. Usip: Unsupervised stable interest point detection from 3d point clouds. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 361–370, 2019. 1, 2

[13] Yongcheng Liu, Bin Fan, Shiming Xiang, and Chunhong Pan. Relation-shape convolutional neural network for point cloud analysis. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8895–8904, 2019. 2

[14] Wan-Chun Ma, Fu-Che Wu, and Ming Ouhyoung. Skeleton extraction of 3d objects with radial basis functions. In *2003 Shape Modeling International.*, pages 207–215. IEEE, 2003. 2

[15] Ajmal S Mian, Mohammed Bennamoun, and Robyn Owens. Three-dimensional model-based object recognition and segmentation in cluttered scenes. *IEEE transactions on pattern analysis and machine intelligence*, 28(10):1584–1601, 2006. 1

[16] 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. 2, 6

[17] Charles R Qi, Li Yi, Hao Su, and Leonidas J Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. In *NIPS*, 2017. 2, 4

[18] Wei Shen, Kai Zhao, Yuan Jiang, Yan Wang, Zhijiang Zhang, and Xiang Bai. Object skeleton extraction in natural images by fusing scale-associated deep side outputs. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 222–230, 2016. 2

[19] Ivan Sipiran and Benjamin Bustos. Harris 3d: a robust extension of the harris operator for interest point detection on 3d meshes. *The Visual Computer*, 27(11):963, 2011. 1, 2, 6

[20] Jian Sun, Maks Ovsjanikov, and Leonidas Guibas. A concise and provably informative multi-scale signature based on heat diffusion. In *Computer graphics forum*, volume 28, pages 1383–1392. Wiley Online Library, 2009. 2

[21] Levente Tamas and Lucian Cosmin Goron. 3d semantic interpretation for robot perception inside office environments. *Engineering Applications of Artificial Intelligence*, 32:76–87, 2014. 1

[22] Leizer Teran and Philippos Mordohai. 3d interest point detection via discriminative learning. In *European Conference on Computer Vision*, pages 159–173, 2014. 6

[23] Federico Tombari, Samuele Salti, and Luigi Di Stefano. Unique signatures of histograms for local surface description. In *European conference on computer vision*, pages 356–369. Springer, 2010. 2

[24] Hanyu Wang, Jianwei Guo, Dong-Ming Yan, Weize Quan, and Xiaopeng Zhang. Learning 3d keypoint descriptors for non-rigid shape matching. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 3–19, 2018. 1

[25] Martin Weinmann. *Reconstruction and analysis of 3D scenes*. Springer, 2016. 1

[26] Wenxuan Wu, Zhongang Qi, and Li Fuxin. Pointconv: Deep convolutional networks on 3d point clouds. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9621–9630, 2019. 2, 6, 8

[27] 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. 1

[28] Yifan Xu, Tianqi Fan, Mingye Xu, Long Zeng, and Yu Qiao. Spidercnn: Deep learning on point sets with parameterized convolutional filters. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 87–102, 2018. 2, 6

[29] Jiaqi Yang, Ke Xian, Yang Xiao, and Zhiguo Cao. Performance evaluation of 3d correspondence grouping algorithms. In *2017 International Conference on 3D Vision (3DV)*, pages 467–476. IEEE, 2017. 1

[30] Li Yi, Vladimir G Kim, Duygu Ceylan, I-Chao Shen, Mengyan Yan, Hao Su, Cewu Lu, Qixing Huang, Alla Sheffer, and Leonidas Guibas. A scalable active framework for region annotation in 3d shape collections. *ACM Transactions on Graphics (ToG)*, 35(6):1–12, 2016. 6

[31] 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. 1, 2, 6, 7

[32] Andy Zeng, Kuan-Ting Yu, Shuran Song, Daniel Suo, Ed Walker, Alberto Rodriguez, and Jianxiong Xiao. Multi-view self-supervised deep learning for 6d pose estimation in the amazon picking challenge. In *2017 IEEE international conference on robotics and automation (ICRA)*, pages 1386–1383. IEEE, 2017. 2

[33] Yu Zhong. Intrinsic shape signatures: A shape descriptor for 3d object recognition. *2009 IEEE 12th International Conference on Computer Vision Workshops, ICCV Workshops*, pages 689–696, 2009. 1, 2, 6