

USIP: Unsupervised Stable Interest Point Detection from 3D Point Clouds

Jiaxin Li* Gim Hee Lee

Department of Computer Science, National University of Singapore

Abstract

In this paper, we propose the USIP detector: an Unsupervised Stable Interest Point detector that can detect highly repeatable and accurately localized keypoints from 3D point clouds under arbitrary transformations without the need for any ground truth training data. Our USIP detector consists of a feature proposal network that learns stable keypoints from input 3D point clouds and their respective transformed pairs from randomly generated transformations. We provide degeneracy analysis of our USIP detector and suggest solutions to prevent it. We encourage high repeatability and accurate localization of the keypoints with a probabilistic chamfer loss that minimizes the distances between the detected keypoints from the training point cloud pairs. Extensive experimental results of repeatability tests on several simulated and real-world 3D point cloud datasets from Lidar, RGB-D and CAD models show that our USIP detector significantly outperforms existing hand-crafted and deep learning-based 3D keypoint detectors. Our code is available at the project website.¹

1. Introduction

3D interest point or keypoint detection refers to the problem of finding stable points with well-defined positions that are highly repeatable on 3D point clouds under arbitrary SE(3) transformations. These detected keypoints play important roles in many computer vision and robotics tasks, where 3D point clouds are widely adopted as the data structure to represent objects and scenes in the 3D space. Examples include geometric registration for 3D object modeling [1] or point cloud-based Simultaneous Localization and Mapping (SLAM) [23], and 3D object [15, 19] or place recognition [34]. In these tasks, the detected keypoints are respectively used as correspondences to compute rigid transformations, and locations to extract representative signatures for efficient retrievals. Hence, a keypoint detector that cannot produce highly repeatable and well-localized

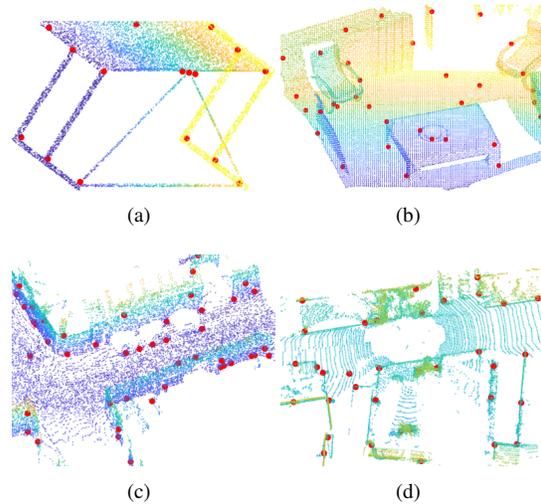


Figure 1. Examples of keypoints detected by our USIP detector on four datasets: (a) ModelNet40 [35], object model. (b) Redwood [5] (Trained on “RGB-D reconstruction dataset” [38]), indoor RGB-D. (c) Oxford RobotCar [21], outdoor SICK LiDAR. (d) KITTI [10] (Trained on Oxford), outdoor Velodyne LiDAR.

keypoints from 3D point clouds under arbitrary transformations would render these tasks to fail catastrophically.

Despite the high number of successful hand-crafted detectors proposed for 2D images [25, 20, 12], significantly lesser hand-crafted detectors [32] with limited success are proposed for hand-crafted detectors on 3D point clouds. This difference can be largely attributed to the difficulty in hand-crafting powerful algorithms to extract meaningful information solely from the Euclidean coordinates of the point cloud in comparison to images that contain richer information from the additional RGB channels. The problem is further aggravated by the fact that it is difficult to hand-craft 3D detectors to handle 3D point clouds in arbitrary transformations, *i.e.*, different reference coordinate frames. In particular, different transformations applied to the same 3D point cloud cause the Euclidean coordinates of each point to change significantly, thus severely affecting the repeatability of the keypoints from the 3D detectors.

It seems evidential that all the above mentioned problems with hand-crafted detectors for 3D point clouds can be resolved by the highly successful data-driven deep net-

*now at nuTonomy: an APTIV company.

¹<https://github.com/lijx10/USIP>

works. However, very few deep learning-based 3D keypoint detectors exist (only one deep learning-based approach [36] exists to date) in contrast to its increasing success on learning 3D keypoint descriptors [7, 6, 38, 16]. This is due to the lack of ground truth training datasets to supervise deep learning-based detectors on 3D point clouds. Unlike 3D descriptors that are supervised by easily available ground truth registered overlapping 3D point clouds [7, 6, 16, 38, 36, 11], it is impossible for anyone to identify and label the “ground truth” keypoints on 3D point clouds. Consequently, most of the works on 3D descriptors [7, 6, 38, 16] ignored the detector problem and are built on top of existing hand-crafted 3D detectors or uniform sampling.

In view of the challenges on both hand-crafted and deep learning-based 3D detectors, we propose the USIP detector: an **U**nsupervised **S**table **I**nterest **P**oint deep learning-based detector that can detect highly repeatable, and accurately localized keypoints from 3D point clouds under arbitrary transformations *without* the need for any ground truth training data. To this end, we design a Feature Proposal Network (FPN) that outputs a set of keypoints and their respective saliency uncertainties from an input 3D point cloud. Our FPN improves keypoint localization by estimating their positions on contrary to existing 3D detectors [29, 36, 39] that select existing points in the point cloud as keypoints, which causes quantization errors. During training, we apply randomly generated SE(3) transformations on each point cloud to get a set of corresponding pairs of transformed point clouds as inputs to the FPN. Furthermore, we identify and prevent the degeneracy of our USIP detector. We encourage high repeatability and accurate localization of the keypoints with a probabilistic chamfer loss that minimizes the distances between the detected keypoints from the training point cloud pairs. Additionally, we introduce a point-to-point loss to enforce the constraint of getting keypoints that lie close to the point cloud. We verify our USIP detector by performing extensive repeatability tests on several simulated and real-world benchmark 3D point cloud datasets from Lidar, RGB-D and CAD models. Some qualitative results are shown in Fig 1.

Our key contributions are summarized as follows:

- Our USIP detector is fully unsupervised, thus avoids the need for ground truth that are impossible to obtain.
- We provide degeneracy analysis of our USIP detector and suggest solutions to prevent it.
- Our FPN improves keypoint localization by estimating the keypoint position instead of choosing it from an existing point in the point cloud.
- We introduce the probabilistic chamfer loss and point-to-point loss to encourage high repeatability and accurate keypoint localization.

- The use of randomly generated transformations on point clouds during training inherently allows our network to achieve good performance under rotations.

2. Related Work

Unlike the recent success of deep learning-based 3D keypoint descriptors [7, 6, 16, 38, 36, 11], most existing 3D keypoint detectors remain hand-crafted. A comprehensive review and evaluation of existing hand-crafted 3D keypoint detectors can be found in [32]. Local Surface Patches (LSP) [3] and Shape Index (SI) [8] are based on the maximum and minimum principal curvatures of a point, and consider the point as a keypoint if it is a global extremum in a pre-defined neighborhood. Intrinsic Shape Signatures (ISS) [39] and KeyPoint Quality (KPG) [22] select salient points that has a local neighborhood with large variations along each principal axis. MeshDoG [37] and Salient Points (SP) [2] construct a scale-space of the curvature with the Difference-of-Gaussian (DoG) operator similar to SIFT [20]. Points with local extrema values over an one-ring neighborhood are selected as keypoints. These methods can be regarded as the 3D extension of SIFT. Laplace-Beltrami Scale-space (LBSS) [33] computes the saliency by applying a Laplace-Beltrami operator on increasing supports for each point.

More recently, LORAX [9] proposes the method of projecting the point set into a depth map and use Principal Component Analysis (PCA) to select keypoints with commonly found geometric characteristics. All hand-crafted approaches share the common trait of relying on the local geometric properties of the points to select keypoints. Hence, the performances of these detectors deteriorate under disturbances such as noise, density variations and/or arbitrary transformations. In contrast, our deep learning-based USIP detector is more resilient to these disturbances by learning from data. To the best of our knowledge, the only existing deep learning-based 3D keypoint detector is the weakly supervised 3DFeatNet [36], which is trained with GPS/INS tagged point clouds. However, the training of 3DFeat-Net is largely focused on learning discriminative descriptors using the Siamese architecture with an attention score map that estimates the saliency of each point as its by-product. It does not ensure good performance of the keypoint detection. In comparison, our USIP is designed to encourage high repeatability and accurate localization of the keypoints. Furthermore, our method is fully unsupervised and does not rely on any form of ground truth datasets.

3. Our USIP Detector

Fig. 2(a) shows the illustration of the pipeline to train our USIP detector. We denote a point cloud from the training dataset as $\mathbf{X} = [X_0, \dots, X_N] \in \mathbb{R}^{3 \times N}$. A set of transformation matrices $\{T_1, \dots, T_L\}$, where $T_l \in$

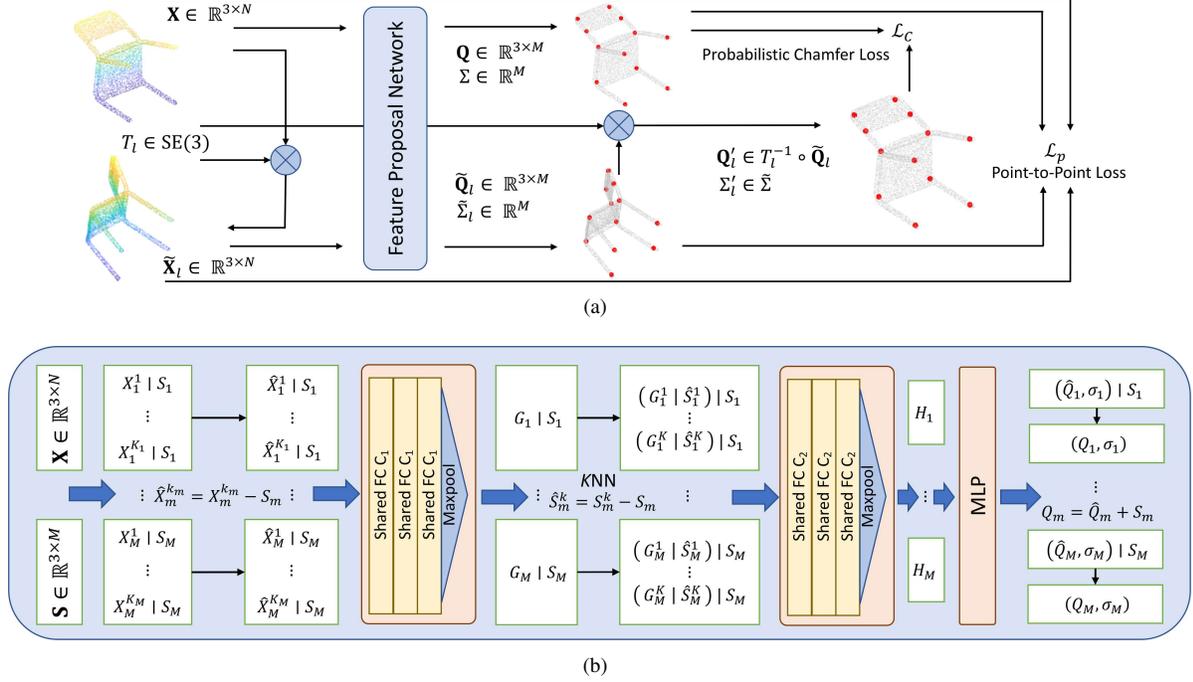


Figure 2. (a) The training pipeline of USIP detector. (b) The architecture of our Feature Proposal Network (FPN). See text for more detail.

$\text{SE}(3)$ is randomly generated and applied to the point cloud \mathbf{X} to form L pairs of training inputs denoted as $\{\{\mathbf{X}, \tilde{\mathbf{X}}_1\}, \dots, \{\mathbf{X}, \tilde{\mathbf{X}}_L\}\}$, where $\tilde{\mathbf{X}}_l = T_l \circ \mathbf{X} \in \mathbb{R}^{3 \times N}$. Here, we use the operator \circ to denote matrix multiplication under homogeneous coordinate with a slight abuse of notation. We drop the indices l for brevity and refer to a triplet of training pair of point clouds and their corresponding transformation matrix as $\{\mathbf{X}, \tilde{\mathbf{X}}, T\}$. During training, \mathbf{X} and $\tilde{\mathbf{X}}$ are respectively fed into the FPN, which outputs M proposal keypoints and its saliency uncertainties denoted as $\{\mathbf{Q} = [Q_1, \dots, Q_M]^T, \Sigma = [\sigma_1, \dots, \sigma_M]^T\}$ and $\{\tilde{\mathbf{Q}} = [\tilde{Q}_1, \dots, \tilde{Q}_M]^T, \tilde{\Sigma} = [\tilde{\sigma}_1, \dots, \tilde{\sigma}_M]^T\}$ for the respective point cloud. $Q_m \in \mathbb{R}^3$, $\tilde{Q}_m \in \mathbb{R}^3$, $\sigma_m \in \mathbb{R}^+$ and $\tilde{\sigma}_m \in \mathbb{R}^+$. We enforce $\sigma_m \in \mathbb{R}^+$ and $\tilde{\sigma}_m \in \mathbb{R}^+$ so that it is a valid rate parameter in our probabilistic chamfer loss (see later paragraph). To improve keypoint localization, it is not necessary for all $Q_m \in \mathbf{Q}$ to be any of the points in \mathbf{X} . Similar condition applies to all $\tilde{Q}_m \in \tilde{\mathbf{Q}}$.

We undo the transformation on $\tilde{\mathbf{Q}}$ with a slight abuse of notation to get $\mathbf{Q}' = T^{-1} \circ \tilde{\mathbf{Q}} \in \mathbb{R}^{3 \times M}$, so that \mathbf{Q}' can be compared directly to \mathbf{Q} . Here, we made an assumption that the saliency uncertainties remain unaffected after the transformation, *i.e.*, $\Sigma' = \tilde{\Sigma}$. The objectives of detecting keypoints that are highly repeatable and accurately localized from 3D point clouds under arbitrary transformations can now be achieved by formulating a loss function that minimizes the difference between \mathbf{Q} and \mathbf{Q}' . To this end, we propose the loss function: $\mathcal{L} = \mathcal{L}_c + \lambda \mathcal{L}_p$, where \mathcal{L}_c is the probabilistic chamfer loss that minimizes the probabilistic distances between all correspondence pairs of keypoints in

\mathbf{Q} and \mathbf{Q}' . \mathcal{L}_p is the point-to-point loss that minimizes the distance of the estimated keypoints to their respective nearest neighbor in the point cloud. This constrains the estimated keypoints to be close to the point cloud. λ is a hyperparameter that adjust the relative contribution of \mathcal{L}_c and \mathcal{L}_p to the total loss. More specifically:

Probabilistic Chamfer Loss A straightforward way to minimize the difference between \mathbf{Q} and \mathbf{Q}' is to use the chamfer loss:

$$\sum_{i=1}^M \min_{Q'_j \in \mathbf{Q}'} \|Q_i - Q'_j\|_2^2 + \sum_{j=1}^M \min_{Q_i \in \mathbf{Q}} \|Q_i - Q'_j\|_2^2, \quad (1)$$

that minimizes the distance of each point in one point cloud with its nearest neighbor in the other point cloud. However, the M proposals are not equally salient. The receptive field of a point Q_i can be a featureless surface since the receptive field is limited to a small volume. In this case, it is detrimental to force the FPN to minimize the distance between Q_i and Q'_j , where Q'_j is the nearest neighbor of Q_i in \mathbf{Q}' .

To mitigate the above problem, we design our FPN to learn the saliency uncertainties Σ and Σ' of the proposal keypoints \mathbf{Q} and \mathbf{Q}' with a probabilistic chamfer loss \mathcal{L}_c . In particular, we propose to formulate \mathcal{L}_c with an exponential distribution that measures the probabilistic distances between \mathbf{Q} and \mathbf{Q}' with the saliency uncertainties Σ and Σ' . More formally, the probability distribution between Q_i and

Q'_j for $i = 1, \dots, M$ is given by:

$$p(d_{ij} | \sigma_{ij}) = \frac{1}{\sigma_{ij}} \exp\left(-\frac{d_{ij}}{\sigma_{ij}}\right), \quad \text{where} \quad (2)$$

$$\sigma_{ij} = \frac{\sigma_i + \sigma'_j}{2} > 0, \quad d_{ij} = \min_{Q'_j \in \mathbf{Q}'} \|Q_i - Q'_j\|_2 \geq 0.$$

$p(d_{ij} | \sigma_{ij})$ is a valid probability distribution since it integrates to 1. A shorter distance d_{ij} between the proposal keypoints Q_i and Q'_j gives a higher probability that Q_i and Q'_j are highly repeatable and accurately localized keypoints in the point clouds \mathbf{X} and $\tilde{\mathbf{X}}$. Assuming i.i.d for all $d_{ij} \in D_{ij}$, the joint distribution between \mathbf{Q} and \mathbf{Q}' is given by:

$$p(D_{ij} | \Sigma_{ij}) = \prod_{i=1}^M p(d_{ij} | \sigma_{ij}). \quad (3)$$

It is important to note that the probability distribution is not symmetrical when the order of the point cloud is swapped, *i.e.*, \mathbf{Q}' and \mathbf{Q} , due to a different set of nearest neighbors, *i.e.*, $d_{ij} \neq d_{ji}$ and $\sigma_{ij} \neq \sigma_{ji}$. Hence, the joint distribution between \mathbf{Q}' and \mathbf{Q} is given by:

$$p(D_{ji} | \Sigma_{ji}) = \prod_{j=1}^M p(d_{ji} | \sigma_{ji}), \quad \text{where} \quad (4)$$

$$\sigma_{ji} = \frac{\sigma'_j + \sigma_i}{2} > 0, \quad d_{ji} = \min_{Q_i \in \mathbf{Q}} \|Q_i - Q'_j\|_2 \geq 0.$$

Finally, the probabilistic chamfer loss \mathcal{L}_c between \mathbf{Q}' and \mathbf{Q} is given by the negative log-likelihood of the joint distributions defined in Eq. 3 and 4:

$$\begin{aligned} \mathcal{L}_c &= \sum_{i=1}^M -\ln p(d_{ij} | \sigma_{ij}) + \sum_{j=1}^M -\ln p(d_{ji} | \sigma_{ji}) \\ &= \sum_{i=1}^M \left(\ln \sigma_{ij} + \frac{d_{ij}}{\sigma_{ij}} \right) + \sum_{j=1}^M \left(\ln \sigma_{ji} + \frac{d_{ji}}{\sigma_{ji}} \right). \end{aligned} \quad (5)$$

We further analyze the physical meaning of σ_{ij} or σ_{ji} by computing the extrema of Eq. 2 from its first derivative over σ_{ij} :

$$\frac{\partial p(d_{ij} | \sigma_{ij})}{\partial \sigma_{ij}} = \frac{d_{ij} \exp(-d_{ij}/\sigma_{ij})}{\sigma_{ij}^3} - \frac{\exp(-d_{ij}/\sigma_{ij})}{\sigma_{ij}^2}, \quad (6)$$

and solve for the stationary points:

$$\frac{\partial p(d_{ij} | \sigma_{ij})}{\partial \sigma_{ij}} = 0 \Rightarrow \sigma_{ij} = d_{ij}. \quad (7)$$

Furthermore, the second derivative $p''(d_{ij} | \sigma_{ij})|_{\sigma_{ij}=d_{ij}} < 0$ means that given a fixed $d_{ij} \neq 0$, the highest probability $p(d_{ij} | \sigma_{ij})$ is achieved at $\sigma_{ij} = d_{ij}$. Consider any triplet

of proposal keypoints $\{Q_i, Q'_j, Q'_k\}$, where d_{ij} and d_{ki} are the distances between the nearest neighbors $\{Q_i, Q'_j\}$ and $\{Q'_k, Q_i\}$ (Q_i can be the nearest neighbor in both orders of \mathbf{Q} and \mathbf{Q}' since chamfer distance is not bijective). σ'_k has to take a large value when $d_{ij} \rightarrow 0$ and d_{kj} is large because we have shown that $\sigma_{ij} = d_{ij}$ and $\sigma_{ki} = d_{ki}$ at optimum. Furthermore, $d_{ij} \rightarrow 0$ and d_{kj} is large implies that $\{Q_i, Q'_j\}$ are repeatable and accurately localized keypoints while Q'_k is not. Hence, a large saliency uncertainty σ'_k for a bad proposal keypoint Q'_k at optimum shows that our probabilistic chamfer loss is guiding the FPN to learn correctly.

Point-to-Point Loss To avoid quantization error in the positions of the keypoints, we design the FPN such that it is not necessary that the proposal keypoints \mathbf{Q} to be any of the points in \mathbf{X} . However, this can cause the FPN to give erroneous proposal keypoints \mathbf{Q} that are far away from the point cloud \mathbf{X} . We circumvent this problem by adding a loss function \mathcal{L}_p that penalizes $Q_m \in \mathbf{Q}$ for being too far from \mathbf{X} . We also apply similar penalty on $\tilde{\mathbf{Q}}$ and $\tilde{\mathbf{X}}$. This loss can be formulated as either the point-to-point loss [1]:

$$\mathcal{L}_{\text{point}} = \sum_{i=1}^M \min_{X_j \in \mathbf{X}} \|Q_i - X_j\|_2^2 + \sum_{i=1}^M \min_{\tilde{X}_j \in \tilde{\mathbf{X}}} \|\tilde{Q}_i - \tilde{X}_j\|_2^2, \quad (8)$$

where $X_j \in \mathbf{X}$ is the nearest neighbor of Q_i or the point-to-plane loss [26, 4]:

$$\mathcal{L}_{\text{plane}} = \sum_{i=1}^M \mathcal{N}_j^T (Q_i - X_j) + \sum_{i=1}^M \tilde{\mathcal{N}}_j^T (\tilde{Q}_i - \tilde{X}_j), \quad (9)$$

where \mathcal{N}_j and $\tilde{\mathcal{N}}_j$ are the nearest surface normal in \mathbf{X} to Q_i and $\tilde{\mathbf{X}}$ to \tilde{Q}_i , respectively. We set $\mathcal{L}_p = \mathcal{L}_{\text{point}}$ by default since we found experimentally that both loss functions give similar performances.

4. Feature Proposal Network

The network architecture of our FPN is shown in Fig. 2(b). We first sample M nodes denoted as $\mathbf{S} = [S_1, \dots, S_M] \in \mathbb{R}^{3 \times M}$ with Farthest Point Sampling (FPS) from a given input point cloud $\mathbf{X} \in \mathbb{R}^{3 \times N}$. A neighborhood of points is built for each node $S_m \in \mathbf{S}$ using point-to-node grouping [18, 17], which is denoted as $\{\{X_1^1 | S_1, \dots, X_1^{K_1} | S_1\}, \dots, \{X_M^1 | S_M, \dots, X_M^{K_M} | S_M\}\}$. K_1, \dots, K_M represents the number of points associated with the each of the nodes in \mathbf{S} . The advantage of point-to-node association over node-to-point k NN search or radius-based ball-search is two-fold: (1) Every point in \mathbf{X} is associated with one node, while some points may be left out in node-to-point k NN search and ball-search. (2) Point-to-node grouping automatically adapts to various scale and

point density, while k NN search and ball-search are vulnerable to density variation and varying scales, respectively. To make FPN *translation equivariant*, we normalize each neighborhood point $\{X_m^1|S_m, \dots, X_m^K|S_m\}$ into $\{\hat{X}_m^1|S_m, \dots, \hat{X}_m^K|S_m\}$ by subtracting from its respective node S_m , i.e., $\hat{X}_m^k = X_m^k - S_m$. Each cluster of normalized local neighborhood points is then fed into a PointNet-like network [24] shown in Fig. 2(b) to get a local feature vector G_m associated with S_m . A k NN grouping layer is applied on the set of local feature vectors $\{G_1|S_1, \dots, G_M|S_M\}$ to achieve hierarchical information aggregation. Specifically, the k nearest neighbors of each pair of $(G_m|S_m)$ are retrieved as $\{(G_m^1|S_m^1)|S_m, \dots, (G_m^K|S_m^K)|S_m\}$. These k NN local feature vectors are then normalized by subtracting with its respective S_m to get a position-independent neighborhood denoted as $\{G_m^1|\hat{S}_m^K|S_m, \dots, G_m^K|\hat{S}_m^K|S_m\}$, where $\hat{S}_m^K = S_m^K - S_m$, before feeding into another network to get a set of feature vectors $\{H_1, \dots, H_M\}$. A simple Multi-Layer Perceptron (MLP) is then used to estimate M proposal keypoints $\{\hat{Q}_1|S_1, \dots, \hat{Q}_M|S_M\}$, where $\hat{Q}_m \in \mathbb{R}^3$, and saliency uncertainties $\{\sigma_1, \dots, \sigma_M\}$, where $\sigma_m \in \mathbb{R}^+$ from $\{H_1, \dots, H_M\}$. Finally, we un-normalize each \hat{Q}_m with S_m , i.e., $Q_m = \hat{Q}_m + S_m$ to get the final proposal keypoints $\{Q_1, \dots, Q_M\}$. It is important to note that the size of the receptive field is controlled by the number of proposals M and K in k NN layers and it determines the level-of-detail for each feature. Large receptive field leads to features that are salient on a large-scale and vice versa.

5. Degeneracy Analysis

Let us denote the FPN as $f(\mathbf{Y}) : \mathbf{Y} \rightarrow \mathbb{R}^{3 \times M}$, where $\mathbf{Y} = [Y_1, \dots, Y_N] \in \mathbb{R}^{3 \times N}$ is the input of the network. We further denote a transformation matrix $T \in \text{SE}(3)$, where $R \in \text{SO}(3)$ and $t \in \mathbb{R}^3$ are the rotation matrix and translation vector in T . We get $\mathbf{Y}' = R\mathbf{Y} \oplus t$, where \oplus is the operator to denote the addition of t to every 3×1 entries of the other term. We say that the network is degenerate when it outputs *trivial solutions* where $f(\mathbf{Y}') \equiv Rf(\mathbf{Y}) \oplus t$ is satisfied for all R and t .

Lemma 1. $f(\mathbf{Y}') \equiv Rf(\mathbf{Y}) \oplus t$ when $f(\cdot)$ outputs the *centroid* of the input point cloud, i.e., $f(\mathbf{Y}) = \frac{1}{N} \sum_n Y_n$ and $f(\mathbf{Y}') = \frac{1}{N} \sum_n Y'_n$.

Proof. Putting $Y'_n = RY_n + t$ into $f(\mathbf{Y}') = \frac{1}{N} \sum_n Y'_n$, we get $f(\mathbf{Y}') = \frac{1}{N} \sum_n (RY_n + t) = R(\frac{1}{N} \sum_n Y_n) + t = Rf(\mathbf{Y}) \oplus t$. Hence, $f(\mathbf{Y}') \equiv Rf(\mathbf{Y}) \oplus t$ which completes our proof that the network degenerates when it outputs the centroid of the input point cloud. \square

Lemma 2. $f(\mathbf{Y}') \equiv Rf(\mathbf{Y}) \oplus t$ when $f(\cdot)$ is *translational equivariant*, i.e., $f(\cdot) \oplus t = f(\cdot \oplus t)$, and outputs points that are in the linear subspace of any *principal axis* from

the input point cloud denoted as $\mathbf{U} = [U_1, U_2, U_3] \in \mathbb{R}^{3 \times 3}$, i.e., $f(\mathbf{Y}) = [c_1 U_i^T, \dots, c_M U_i^T]^T$ and

$$\begin{aligned} f(\mathbf{Y}') &= f(R\mathbf{Y} \oplus t) \\ &= f(R\mathbf{Y}) \oplus t \quad (\text{translation equivariance}) \quad (10) \\ &= [c_1 U_i'^T, \dots, c_M U_i'^T]^T \oplus t, \end{aligned}$$

where U_i can be any principal axis in \mathbf{U} and c_1, \dots, c_M are scalar coefficients in \mathbb{R} .

Proof. Let $V = \frac{1}{N} \sum_n (Y_n - \bar{Y})(Y_n - \bar{Y})^T$ and $V' = \frac{1}{N} \sum_n (Y'_n - \bar{Y}')(Y'_n - \bar{Y}')^T$ denote the covariance matrices of \mathbf{Y} and \mathbf{Y}' , respectively. $\bar{Y} = \frac{1}{N} \sum_n Y_n$ and $\bar{Y}' = \frac{1}{N} \sum_n Y'_n$ are the centroids of \mathbf{Y} and \mathbf{Y}' , respectively. Putting $Y'_n = RY_n + t$ into \bar{Y}' and V' , we get:

$$V' = R \frac{1}{N} \sum_n (Y_n - \bar{Y})(Y_n - \bar{Y})^T R^T = RV R^T. \quad (11)$$

Taking the Singular Value Decomposition (SVD) of V and V' , we get $V = \mathbf{U}\mathbf{D}\mathbf{U}^T$ and $V' = \mathbf{U}'\mathbf{D}'\mathbf{U}'^T$, where \mathbf{D} and \mathbf{D}' are the 3×3 diagonal matrices of singular values, and \mathbf{U} and \mathbf{U}' are the 3×3 Eigenvectors that are also the principal axes of \mathbf{Y} and \mathbf{Y}' , respectively. Putting the SVD of V and V' into Eq. 11, we get:

$$\begin{aligned} V' &= RV R^T = R\mathbf{U}\mathbf{D}\mathbf{U}^T R^T = (R\mathbf{U})\mathbf{D}(R\mathbf{U})^T \\ &\equiv \mathbf{U}'\mathbf{D}'\mathbf{U}'^T \Rightarrow \mathbf{U}' = R\mathbf{U}. \end{aligned} \quad (12)$$

Putting the relationship from Eq. 12 into $f(\mathbf{Y}') = [c_1 U_i'^T, \dots, c_M U_i'^T]^T \oplus t$, we get:

$$f(\mathbf{Y}') = R[c_1 U_i^T, \dots, c_M U_i^T]^T \oplus t \equiv Rf(\mathbf{Y}) \oplus t, \quad (13)$$

which completes our proof that the network degenerates when it outputs a set of points on any principal axis. \square

Discussions We note that the network requires sufficient global semantic information of the input point cloud, e.g., the input is the whole point cloud or clusters of local neighbor points that contain large receptive fields, to learn the trivial solutions of centroid or set of points on the principal axes. Hence, the degeneracies can be easily prevented by limiting the receptive fields of the FPN. We achieve this by setting the number of clusters M and K nearest neighbors of the clusters in the FPN (refer to Sec. 4 for the definitions of M and K) to reasonable values. Small values for M or high values for K increases the receptive field and causes the FPN to degenerate. Fig. 3 show some examples of the degeneracies with different K values at $M = 64$. It is interesting to note that the principal axis degeneracy occurs when K is set to a mid-range value, and centroid degeneracy occurs when K is set to a high value. This implies that larger receptive fields, i.e., a higher global semantic information is needed for the network to learn the centroid. We also notice experimentally that the degeneracies

(both centroid and principal axis) occur in point clouds with more regular shapes, *e.g.* objects from ModelNet40 where the centroid and principal axes are more well-defined.

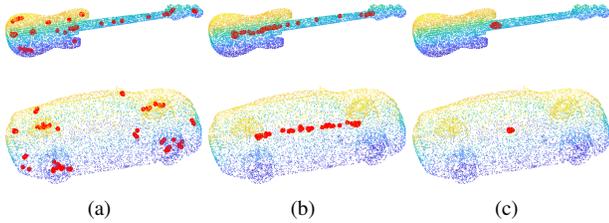


Figure 3. Increasing K values in FPN causes degeneracies ($M = 64$). (a) No degeneracy with $K = 9$ (low value). (b) Principal axis degeneracy with $K = 24$ (mid-range value). (c) Centroid degeneracy with $K = 64$ (high value).

6. Experiments

Following [32], we evaluate the *repeatability* (Sec. 6.1), *distinctiveness* (Sec. 6.2) and *computational efficiency* (Sec. 6.3) of our USIP detector on 4 datasets from object models, outdoor Lidar and indoor RGB-D scans. Additionally, we compare our evaluations to existing detectors - ISS [39], Harris-3D [12], SIFT-3D [20] and 3DFeat-Net [36].

Implementation Details Three USIP detectors are respectively trained for outdoor Lidars, RGB-D scans and object models. Specifically, we use the Oxford [21] for outdoor Lidar, “RGB-D reconstruction dataset” [38] for RGB-D, and ModelNet40 [35] for object models. The PCL [29] implementations of the classical detectors, *i.e.*, ISS, Harris-3D and SIFT-3D are used for the comparisons. We take the pretrained models of 3DFeat-Net [36] for KITTI [10] and Oxford, and train separate models for Redwood and ModelNet40 using the codes provided by 3DFeat-Net.

Qualitative Visualization Fig. 7 shows some results from our USIP detector on ModelNet40. We can clearly see that our USIP learns keypoints on corners, edges, center of small surfaces, etc. Keypoints in the first row of Fig. 7 are selected with Non-Maximum Suppression (NMS) and thresholding on the saliency uncertainty σ . In the second row, keypoints are selected with only NMS. Keypoints with small σ are shown in bright red and get darker with larger σ .

6.1. Repeatability

Repeatability refers to the ability of a detector to detect keypoints in the same locations under various disturbances such as view-point variations, noise, missing parts, etc. It is often taken as the most important measure of keypoint detectors because it is a standalone measure that depends only on the detector (without a descriptor). Given two point clouds $\{\mathbf{X}, \tilde{\mathbf{X}}\}$ of a scene captured from different view-points such that $\{\mathbf{X}, \tilde{\mathbf{X}}\}$ are related by a rotation matrix

	KITTI	Oxford	Redwood	ModelNet40
Type	Velodyne lidar	SICK lidar	RGB-D	CAD Model
Scale	200m	60m	10m	2
# point	16,384	16,384	10,240	5,000
ϵ in Eq. 14	0.5m	0.5m	0.1m	0.03
Rotation	2D	2D	3D	3D
Noise	Sensor	Sensor	Gaussian	Gaussian
Occlusion	Yes	Yes	Yes	No
Density Variation	Yes	No	No	No
Missing Parts	Yes	Yes	Yes	No

Table 1. Datasets used in evaluating keypoint repeatability.

$R \in \text{SO}(3)$ and a translational vector $t \in \mathbb{R}^3$. A keypoint detector detects a set of keypoints $\mathbf{Q} = [Q_1, \dots, Q_M]$ and $\tilde{\mathbf{Q}} = [\tilde{Q}_1, \dots, \tilde{Q}_M]$ from $\{\mathbf{X}, \tilde{\mathbf{X}}\}$, respectively. A keypoint $Q_i \in \mathbf{Q}$ is repeatable if the distance between $RQ_i + t$ and its nearest neighbor $\tilde{Q}_j \in \tilde{\mathbf{Q}}$ is less than a threshold ϵ , *i.e.*,

$$\|RQ_i + t - \tilde{Q}_j\|_2 < \epsilon. \quad (14)$$

Test Datasets We evaluate repeatability on four test datasets - KITTI, Oxford, Redwood and ModelNet40. Note that our USIP is not trained on KITTI nor Redwood. We use the KITTI and Oxford test datasets prepared by 3DFeat-Net [36]. Each pair of point clouds $\{\mathbf{X}, \tilde{\mathbf{X}}\}$ are captured from nearby locations of within 10m and manually augmented with random 2D rotations. $\{\mathbf{X}, \tilde{\mathbf{X}}\}$ in Redwood are from simulated RGB-D cameras with 3D rotations / translations and Gaussian noise. The overlapped areas between $\{\mathbf{X}, \tilde{\mathbf{X}}\}$ are as low as 30%. In ModelNet40, $\tilde{\mathbf{X}}$ is obtained by augmenting \mathbf{X} with random 3D rotations. Points in KITTI, Oxford and Redwood are in its original scale while points in ModelNet40 are normalized to $[-1, 1]$. Details of the datasets are shown in Tab. 1. The scale refers to the diameter of the point clouds.

Relative Repeatability We use relative repeatability that normalizes over the total number of detected keypoints $|\mathbf{Q}|$ for fair comparisons, *i.e.*, $\text{repeatability} = |\mathbf{Q}_{\text{rep}}|/|\mathbf{Q}|$, where \mathbf{Q}_{rep} is the number of keypoints that passed the repeatability test in Eq. 14. We set the parameters of each keypoint detector in each dataset to generate 4, 8, 16, 32, 64, 128, 256 and 512 keypoints or close to these numbers when it is not possible to set the detectors (SIFT-3D, Harris-3D and ISS) to generate exact number of keypoints. Note that in general the repeatability should be proportional to the number of keypoints. In the extreme case that $\mathbf{Q} = \mathbf{X}$, *i.e.*, each point is regarded as a keypoint, the repeatability is the same as the percentage of overlap between $\{\mathbf{X}, \tilde{\mathbf{X}}\}$. As shown in Fig. 4, our USIP generally outperforms other detectors by a significant margin on the 4 datasets over 8 different number of keypoints. In the extremely hard case that only 4 keypoints are detected, our method achieved relative repeatability of 34%, 23%, 10% and 60% for KITTI, Oxford, Redwood and ModelNet40, respectively. In the case of 64

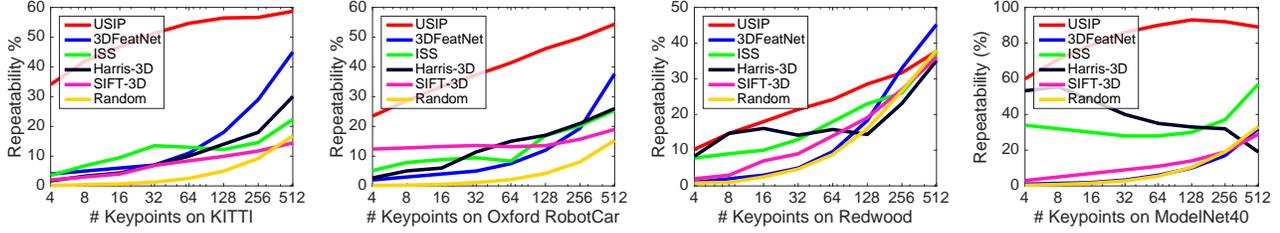


Figure 4. Relative repeatability when different number of keypoints are detected. Left to right: KITTI, Oxford, Redwood, ModelNet40.

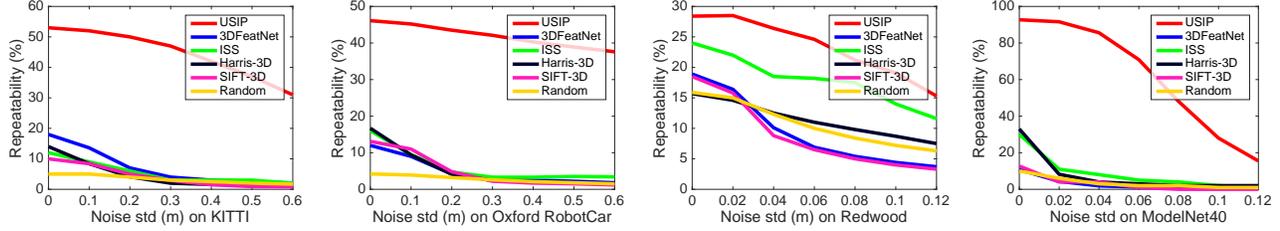


Figure 5. Relative repeatability when Gaussian noise $\mathcal{N}(0, \sigma_{noise})$ is added to the input point clouds. Keypoint number is fixed to 128.

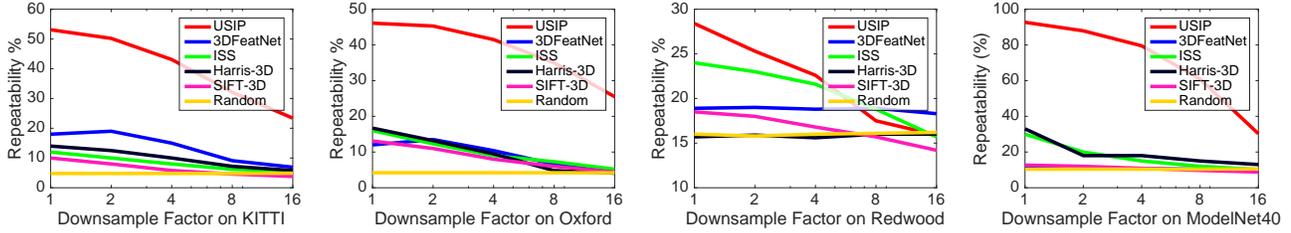


Figure 6. Relative repeatability when the input point cloud is randomly downsampled by some factors. Keypoint number is fixed to 128.

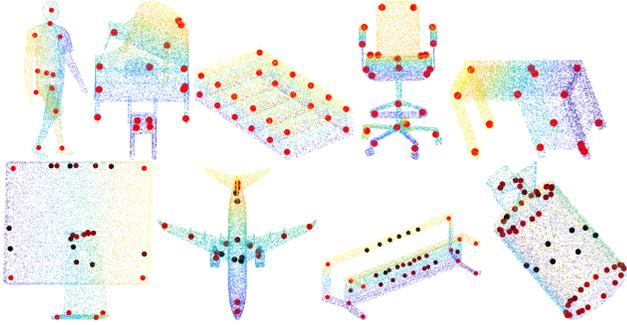


Figure 7. Examples of keypoints from our USIP on ModelNet40.

keypoints, our performance is roughly 4.2x, 2.8x, 1.3x and 2.6x higher than the second best detector.

Robustness to Noise The original points in KITTI and Oxford are already corrupted with sensor noise. We further augment the point clouds in the 4 datasets with Gaussian noise $\mathcal{N}(0, \sigma_{noise})$, where σ_{noise} is up to 0.6m for KITTI and Oxford, 0.12m for Redwood and 0.12 (no unit) for ModelNet40. The number of keypoints is fixed to 128. Our USIP is a lot more robust than other detectors as shown in Fig. 5. In KITTI and Oxford, the performances of other detectors fall to the level of random sampling when $\sigma_{noise} \geq 0.2m$, while our USIP does not show significant drop in performance even with $\sigma_{noise} \geq 0.6m$. In Redwood, other

methods except USIP and ISS deteriorate to random sampling with $\sigma_{noise} \geq 0.02m$. In ModelNet40, our method maintain high repeatability of 91% with $\sigma_{noise} = 0.02$, while all other methods drop below 8%.

Robustness to Downsampling We evaluate the repeatability of the detectors on input point clouds downsampled by some factors using random selection. The results are shown in Fig. 6, where the down-sample factor denoted as α means the number of points is reduced to $\frac{1}{\alpha}$ of the original number shown in Tab. 1. We can see that the repeatability of our USIP remains satisfactory even with a $16\times$ down-sampling on KITTI, Oxford and ModelNet40. The only exception is the Redwood dataset, where almost all detectors perform poorly on high downsample factors. Indoor RGB-D scans in Redwood consist of many large and flat surface like wall, ceiling, etc. Furthermore, there are very few distinguishable and non-occluded structures, which are further aggravated by severe downsampling. Hence, it is difficult to detect repeatable keypoints with these RGB-D scans.

6.2. Distinctiveness: Point Cloud Registration

Distinctiveness is a measure of the performance of keypoint detectors and descriptors for finding correspondences in point cloud registration. Hence, distinctiveness is not as good as repeatability as an evaluation criterion on keypoint detectors because it is confounded with the performance of

	Registration Failure Rate (%)				Inlier Ratio (%)			
	Our Desc.	3DFeatNet[36]	FPFH[27]	SHOT[31]	Our Desc.	3DFeatNet	FPFH	SHOT
Random	18.83	42.14	49.95	68.39	7.47	4.48	5.45	4.46
SIFT-3D[29, 20]	15.44	42.63	79.72	84.49	7.36	5.47	4.24	4.11
ISS[29, 39]	5.97	25.96	37.09	69.83	8.52	4.71	4.44	3.45
Harris-3D[29, 12]	3.81	13.56	49.49	51.29	10.57	6.58	4.78	5.00
3DFeatNet[36]	2.61	2.26	12.15	11.76	15.66	10.76	9.55	8.46
USIP	1.41	1.55	8.37	5.40	32.20	22.48	18.77	18.21

Table 2. Point cloud registration results on KITTI. The number of keypoints is fixed to 256.

the descriptor. We mitigate this limitation by evaluating point cloud registration over several existing keypoint descriptors. We also use the results to show that our USIP detector works with different existing keypoint descriptors.

Experiment Setup We follow the point cloud registration pipeline from 3DFeat-Net [36] on their KITTI test dataset. Four descriptors are used to perform keypoint description, *i.e.*, three off-the-shelf descriptors: 3DFeatNet, FPFH [28], SHOT[31], and our own descriptor inspired by 3DFeat-Net with minor modifications, which is denoted as “Our Desc.” (details are in our supplementary material). Registration of a pair of point clouds involves 4 steps: (a) Extract keypoints and their corresponding descriptor vectors from each point cloud. (b) Establish keypoint-to-keypoint correspondences by nearest neighbor search of the descriptor vectors. (c) Perform RANSAC on the two matched keypoint sets to find the rotation and translation that have the most inliers. (d) Compare the resulted rotation and translation with the ground truth. A pair of point cloud is regarded as successfully registered if Relative Translational Error (RTE) $< 2m$, and Relative Rotation Error (RRE) $< 5^\circ$.

Registration Results We perform registration evaluations over the combination of 6 keypoint detectors and 4 descriptors. The registration failure rate and keypoint inlier ratio are shown in Tab. 2. Compared to other detectors, our USIP achieves the lowest registration failure rate and the highest inlier ratio with a considerable margin on all the 4 descriptors. The significance of the results in Tab. 2 is two fold. First, our USIP works well with various hand-crafted (FPFH and SHOT) and deep learning-based (our desc. and 3DFeat-Net) descriptors. Second, our USIP produces more distinctive keypoints since it consistently outperforms other keypoint detectors over the different descriptors on registration failure rate and keypoint inlier ratio as shown in Tab. 2. The experimental configurations in Tab. 2 is not the optimal setting for our USIP detector and descriptor nor the 3DFeat-Net because we have to fix the number of keypoints for fair comparison. In Tab. 3, we illustrate the best registration results for our USIP and 3DFeatNet on KITTI without limitation on the number of keypoints. We again achieve lower failure rate and higher inlier ratio. In addition, we show the visualization of keypoint matching results of two examples

from KITTI and Oxford in Fig. 8.

Detector	Descriptor	Fail(%)	Inlier(%)	RTE(m)	RRE ($^\circ$)
3DFeat-Net	3DFeat-Net	0.57	12.9	0.26 ± 0.26	0.56 ± 0.46
USIP	Our Desc.	0.24	28.0	0.21 ± 0.24	0.42 ± 0.32

Table 3. Point cloud registration on KITTI from the optimal configurations of 3DFeat-Net and our USIP.

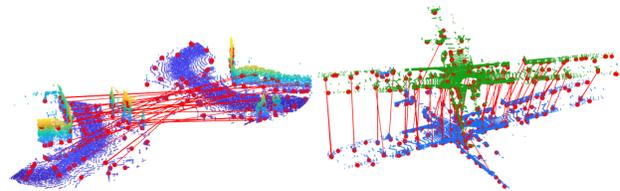


Figure 8. Keypoints and matches from our USIP detector and “Our Desc.”. Best view with color and zoom-in.

6.3. Computational Efficiency

Hand-crafted detectors are deployed with single thread C++ codes on an Intel i7 6950X CPU. Our USIP and 3DFeatNet are deployed on a Nvidia 1080Ti, with PyTorch and TensorFlow, respectively. Computational efficiency is evaluated on 2,391 KITTI point clouds, where each point cloud is downsampled to 16,384 points. We record the average time taken to extract 128 keypoints from each point cloud. As shown in Tab. 4, our USIP is an order of magnitude faster than other detectors except random sampling.

Random	SIFT-3D	ISS	Harris-3D	3DFeatNet	USIP
0.0005	0.163	0.388	0.150	0.438	0.011

Table 4. Average time (in seconds) to extract 128 keypoints from KITTI point clouds respectively downsampled to 16,384 points.

7. Conclusion

In this paper, we present the USIP detector, an unsupervised deep learning-based keypoint detector for 3D point clouds. A probabilistic chamfer loss is proposed to guide the network to learn highly repeatable keypoints. We provide mathematical analysis and solutions for network degeneracy, which are supported by experimental results. Extensive evaluations are performed with Lidar scans, RGB-D images and CAD models. Our USIP detector outperforms existing detectors by a significant margin in terms of repeatability, distinctiveness and computational efficiency.

References

- [1] P. J. Besl and N. D. McKay. Method for registration of 3-d shapes. In *Sensor Fusion IV: Control Paradigms and Data Structures*, volume 1611, pages 586–607. International Society for Optics and Photonics, 1992. **1, 4**
- [2] U. Castellani, M. Cristani, S. Fantoni, and V. Murino. Sparse points matching by combining 3d mesh saliency with statistical descriptors. In *Computer Graphics Forum*, volume 27, pages 643–652. Wiley Online Library, 2008. **2**
- [3] H. Chen and B. Bhanu. 3d free-form object recognition in range images using local surface patches. *Pattern Recognition Letters*, 28(10):1252–1262, 2007. **2**
- [4] Y. Chen and G. Medioni. Object modelling by registration of multiple range images. *Image and vision computing*, 10(3):145–155, 1992. **4**
- [5] S. Choi, Q.-Y. Zhou, and V. Koltun. Robust reconstruction of indoor scenes. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5556–5565, 2015. **1**
- [6] H. Deng, T. Birdal, and S. Ilic. Ppf-foldnet: Unsupervised learning of rotation invariant 3d local descriptors. *arXiv preprint arXiv:1808.10322*, 2018. **2**
- [7] H. Deng, T. Birdal, and S. Ilic. Ppfnet: Global context aware local features for robust 3d point matching. *Computer Vision and Pattern Recognition (CVPR), IEEE*, 1, 2018. **2**
- [8] C. Dorai and A. K. Jain. Cosmos-a representation scheme for 3d free-form objects. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 19(10):1115–1130, 1997. **2**
- [9] G. Elbaz, T. Avraham, and A. Fischer. 3d point cloud registration for localization using a deep neural network auto-encoder. In *Computer Vision and Pattern Recognition (CVPR), 2017 IEEE Conference on*, pages 2472–2481. IEEE, 2017. **2**
- [10] A. Geiger, P. Lenz, and R. Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2012. **1, 6**
- [11] Z. Gojcic, C. Zhou, J. D. Wegner, and A. Wieser. The perfect match: 3d point cloud matching with smoothed densities. *arXiv preprint arXiv:1811.06879*, 2018. **2**
- [12] C. G. Harris, M. Stephens, et al. A combined corner and edge detector. In *Alvey vision conference*, volume 15, pages 10–5244. Citeseer, 1988. **1, 6, 8**
- [13] B.-S. Hua, Q.-H. Pham, D. T. Nguyen, M.-K. Tran, L.-F. Yu, and S.-K. Yeung. Scenenn: A scene meshes dataset with annotations. In *3D Vision (3DV), 2016 Fourth International Conference on*, pages 92–101. IEEE, 2016. **12**
- [14] A. E. Johnson and M. Hebert. Using spin images for efficient object recognition in cluttered 3d scenes. *IEEE Transactions on Pattern Analysis & Machine Intelligence*, (5):433–449, 1999. **13**
- [15] R. V. J.W. Tangelder. A survey of content based 3d shape retrieval methods. 2008. **1**
- [16] M. Khoury, Q.-Y. Zhou, and V. Koltun. Learning compact geometric features. In *Proc. of the IEEE Conf. on Computer Vision and Pattern Recognition*, pages 153–61, 2017. **2, 11, 13**
- [17] T. Kohonen. The self-organizing map. *Neurocomputing*, 21(1):1–6, 1998. **4**
- [18] J. Li, B. M. Chen, and G. H. Lee. So-net: Self-organizing network for point cloud analysis. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9397–9406, 2018. **4**
- [19] Z. Lian and A. A. Godil. A comparison of methods for non-rigid 3d shape retrieval. 2012. **1**
- [20] D. G. Lowe. Distinctive image features from scale-invariant keypoints. *International journal of computer vision*, 60(2):91–110, 2004. **1, 2, 6, 8**
- [21] W. Maddern, G. Pascoe, C. Linegar, and P. Newman. 1 Year, 1000km: The Oxford RobotCar Dataset. *The International Journal of Robotics Research (IJRR)*, 36(1):3–15, 2017. **1, 6**
- [22] A. Mian, M. Bennamoun, and R. Owens. On the repeatability and quality of keypoints for local feature-based 3d object retrieval from cluttered scenes. *International Journal of Computer Vision*, 89(2-3):348–361, 2010. **2**
- [23] M. Montemerlo, S. Thrun, D. Koller, B. Wegbreit, et al. Fastslam: A factored solution to the simultaneous localization and mapping problem. *Aaai/iaai*, 593598, 2002. **1**
- [24] C. R. Qi, H. Su, K. Mo, and L. J. Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. *arXiv preprint arXiv:1612.00593*, 2016. **5**
- [25] E. Rublee, V. Rabaud, K. Konolige, and G. Bradski. Orb: An efficient alternative to sift or surf. In *Computer Vision (ICCV), 2011 IEEE international conference on*, pages 2564–2571. IEEE, 2011. **1**
- [26] S. Rusinkiewicz and M. Levoy. Efficient variants of the icp algorithm. In *3-D Digital Imaging and Modeling, 2001. Proceedings. Third International Conference on*, pages 145–152. IEEE, 2001. **4**
- [27] R. B. Rusu, N. Blodow, and M. Beetz. Fast point feature histograms (fpfh) for 3d registration. In *Robotics and Automation, 2009. ICRA'09. IEEE International Conference on*, pages 3212–3217. Citeseer, 2009. **8, 13**
- [28] R. B. Rusu, G. Bradski, R. Thibaux, and J. Hsu. Fast 3d recognition and pose using the viewpoint feature histogram. In *Intelligent Robots and Systems (IROS), 2010 IEEE/RSJ International Conference on*, pages 2155–2162. IEEE, 2010. **8**
- [29] R. B. Rusu and S. Cousins. 3d is here: Point cloud library (pcl). In *2011 IEEE International Conference on Robotics and Automation*, pages 1–4, May 2011. **2, 6, 8**
- [30] F. Tombari, S. Salti, and L. Di Stefano. Unique shape context for 3d data description. In *Proceedings of the ACM workshop on 3D object retrieval*, pages 57–62. ACM, 2010. **13**
- [31] F. Tombari, S. Salti, and L. Di Stefano. Unique signatures of histograms for local surface description. In *European conference on computer vision*, pages 356–369. Springer, 2010. **8**
- [32] F. Tombari, S. Salti, and L. Di Stefano. Performance evaluation of 3d keypoint detectors. *International Journal of Computer Vision*, 102(1-3):198–220, 2013. **1, 2, 6**
- [33] R. Unnikrishnan and M. Hebert. Multi-scale interest regions from unorganized point clouds. In *2008 IEEE Computer Society Conference on Computer Vision and Pattern Recognition Workshops*, pages 1–8. IEEE, 2008. **2**

- [34] M. A. Uy and G. H. Lee. Pointnetvlad: Deep point cloud based retrieval for large-scale place recognition. 2018. [1](#)
- [35] Z. Wu, S. Song, A. Khosla, F. Yu, L. Zhang, X. Tang, and J. 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](#), [6](#)
- [36] Z. J. Yew and G. H. Lee. 3dfeat-net: Weakly supervised local 3d features for point cloud registration. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 607–623, 2018. [2](#), [6](#), [8](#), [11](#), [12](#), [13](#)
- [37] A. Zaharescu, E. Boyer, K. Varanasi, and R. Horaud. Surface feature detection and description with applications to mesh matching. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 373–380. IEEE, 2009. [2](#)
- [38] A. Zeng, S. Song, M. Nießner, M. Fisher, J. Xiao, and T. Funkhouser. 3dmatch: Learning local geometric descriptors from rgb-d reconstructions. In *Computer Vision and Pattern Recognition (CVPR), 2017 IEEE Conference on*, pages 199–208. IEEE, 2017. [1](#), [2](#), [6](#), [11](#), [12](#), [13](#), [19](#)
- [39] Y. Zhong. Intrinsic shape signatures: A shape descriptor for 3d object recognition. In *Computer Vision Workshops (ICCV Workshops), 2009 IEEE 12th International Conference on*, pages 689–696. IEEE, 2009. [2](#), [6](#), [8](#), [13](#)

A. Overview

We provide more details on the algorithms and experiments described in the main paper. Sec. B presents more examples of the network degeneracy. Sec. C evaluates the effect of point-to-point loss \mathcal{L}_p on the keypoint repeatability. Sec. D illustrates the details of our feature descriptor design. Sec. E gives more experiments on point cloud registration tasks. Sec. F presents visualizations of our USIP keypoints in various datasets.

B. More Examples on Degeneracy

As analyzed in Sec. 5, our FPN degenerates when the receptive field becomes sufficiently large, *i.e.*, it has gained sufficient global semantic information. The receptive field of the FPN is controlled by two parameters: number of keypoint proposals M and number of neighbors K in the K NN feature aggregation. More specifically, the receptive field size is proportional to K and inversely proportional to M . In this section, we visualize the network degeneracy by gradually enlarging the receptive field. Fig. 14 shows the degeneracies when $M = 64$ and $K = \{9, 24, 32, 40, 48, 64\}$. Fig. 15 shows the degeneracies when $K = 9$ and $M = \{64, 24, 20, 16, 12, 9\}$.

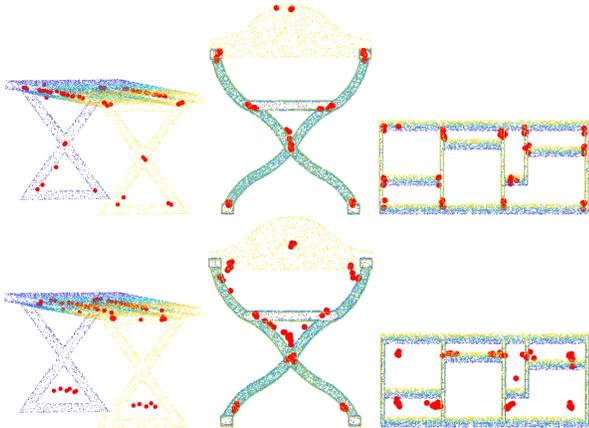


Figure 9. Visualization of USIP keypoints with different λ in Point-to-Point loss. First row $\lambda = 6$, second row $\lambda = 0$.

C. Effect of λ in Point-to-Point Loss \mathcal{L}_p

Sec. 3 of the main paper describes the point-to-point loss \mathcal{L}_p to penalize $Q_m \in \mathbf{Q}$ for being too far from \mathbf{X} . The point-to-point loss \mathcal{L}_p is added to the loss function with the weight λ . Here, we show that our USIP is very robust to the value of λ . Specifically, the repeatability of our USIP keypoints remains almost the same over a wide range of values for λ . Keypoint repeatability is illustrated in Fig. 10 with various λ . Fig. 10 shows that the USIP keypoints are highly repeatable even when λ is small. This is probably because our design to limit the receptive field already guides

the network to learn repeatable keypoints even without the point-to-point loss. On the other hand, the network fails to converge when λ is too large because the point-to-point loss dominates the training process. Nonetheless, training the network without the point-to-point loss does not ensure the keypoints to be close to the input point cloud. The top row of Fig. 9 shows keypoints from our USIP detector trained with $\lambda = 6$, *i.e.*, with point-to-point loss. They are close to the input point cloud. In comparison, the bottom row of Fig. 9 shows from our USIP detector trained without point-to-point loss, *i.e.*, $\lambda = 0$. These are less desirable keypoints that are farther from the input point cloud.

D. Our Descriptor a.k.a “Our Desc.”

Fig 11 shows the network design of “Our Desc.” inspired by 3DFeat-Net [36] as mentioned in Sec. 6.2 of the main paper. Given the output (\mathbf{Q}, Σ) from FPN, a ball $\Omega_m(Q_m, r)$ of points from the point cloud \mathbf{X} within a radius r is built around each $Q_m \in \mathbf{Q}$. A keypoint descriptor $f_m \in \mathbb{R}^L$ is extracted for each Ω_m . The descriptor can be trained with either weak [36] or strong supervision [38, 16]. We improve the keypoint descriptor training by utilizing the keypoint saliency uncertainty Σ in Sec. D.1, D.2, and E.

D.1. Weak Supervision

Weak supervision of the descriptor is based on a triplet loss and the ground truth coarse registrations of the point clouds in the training dataset. Similar to [36], point clouds from the dataset are selected as the anchor samples during training. All overlapping pairs of point clouds to the anchor are defined as positive samples, while non-overlapping pairs of point clouds are defined as the negative samples. We denote the sets of keypoint descriptors extracted from the anchor, positive and negative samples as F_{anc} , F_{pos} and F_{neg} , respectively. We generate these training samples from the Oxford RobotCar and KITTI datasets. More formally, the triplet loss is given by:

$$\mathcal{L}_{\text{dc}}^w = \sum_{m=1}^M w_m \left[\min_{f_i \in F_{\text{pos}}} \|f_m - f_i\|_2 - \min_{f_j \in F_{\text{neg}}} \|f_m - f_j\|_2 + \gamma \right]_+, \quad (15)$$

where $f_m \in F_{\text{anc}}$ is a descriptor from the anchor sample. For each descriptor $f_m \in F_{\text{anc}}$, we minimize the Euclidean distance to its nearest neighbor $f_i \in F_{\text{pos}}$ and maximize the Euclidean distance to its nearest neighbor $f_j \in F_{\text{neg}}$. In addition, a normalized weight w_m is added to our triplet loss. w_m is derived from our USIP keypoint saliency uncertainty σ_m that indicates the reliability of Q_m and f_m . More specifically:

$$w_m = M \cdot \frac{\hat{w}_m}{\sum_{j=1}^M \hat{w}_j}, \quad \hat{w}_m = [\xi - \sigma_m]_+, \quad (16)$$

where ξ is a threshold serves as the upper bound of σ_m .

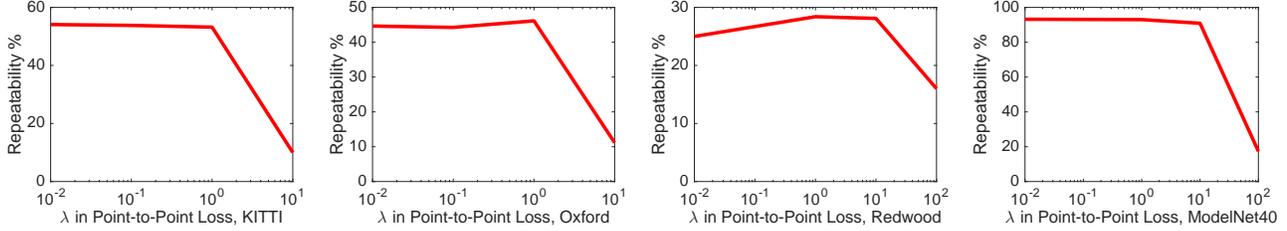


Figure 10. Relative repeatability with different weight λ for the Point-to-Point Loss \mathcal{L}_p . Number of keypoints is fixed to 128. Left to right: KITTI, Oxford, Redwood, ModelNet40.

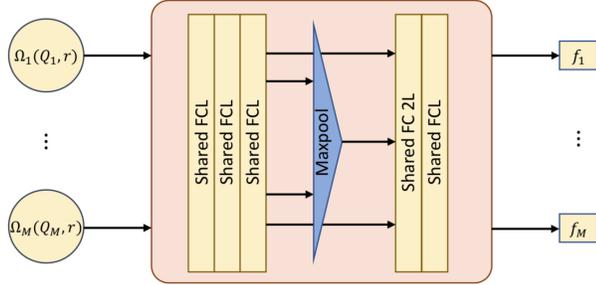


Figure 11. Network architecture of “Our Desc.”.

D.2. Strong Supervision

We do strong supervision of the descriptor network on datasets with ground truth poses, *i.e.*, SceneNN [13] and “3D reconstruction dataset” [38]. The loss function for strong supervision defined on a pair of overlapping point clouds \mathbf{X} and \mathbf{X}' with ground truth poses $G \in \text{SE}(3)$ and $G' \in \text{SE}(3)$ is given by:

$$\mathcal{L}_{\text{dc}}^s = \sum_{m=1}^M w_m \left[\|f_m - f'_i\|_2 - \|f_m - f'_j\|_2 + \gamma \right]_+ \quad (17)$$

f_m and f'_i are keypoint descriptors from \mathbf{X} and \mathbf{X}' , respectively. Additionally, f'_j is a descriptor with keypoint location Q'_j that is within a distance ρ from the keypoint location Q_m of the descriptor f_m , *i.e.*, $\|Q_m - GG'^{-1}Q'_j\|_2 < \rho$. To achieve hard negative mining, we randomly select 50% of f'_j from \mathbf{X}' with the distance between the keypoint locations Q'_j and Q_m larger than ρ . The other 50% are chosen from keypoints with shortest but larger than ρ keypoint distances to Q_m .

E. More Point Cloud Registration Results

We follow the experimental setup and pipeline of 3DFeat-Net [36] to provide more evaluation results on point cloud registration. More specifically, we compare the performance of our USIP detector and “Our Desc.” with other existing keypoint detector and descriptors. The evaluations are done on the Oxford RobotCar and KITTI datasets prepared by [36]. Refer to Sec. 6.2 of the main paper for the

details of the registration steps. A fixed number of 256 keypoints is extracted from each point cloud. We extract the keypoints without Non-Maximum-Suppression (NMS). Furthermore, keypoints with high saliency uncertainty, *i.e.*, large σ , are filtered out.

Datasets The Oxford RobotCar consists of 40 traversals on the same route over a year. 3D point clouds are built by accumulating the 2D scans from SICK LMS-151 LiDAR with the GPS/INS readings. We use 35 traversals, *i.e.*, 21,875 point clouds for training. The remaining 5 traversals, *i.e.*, 828 point clouds and 3,426 overlapping pairs are used for evaluation. Random rotations around the up-axis are applied to each evaluation point cloud. In KITTI, 3D point clouds are directly provided by a Velodyne HDL-64E. We use the 2,831 overlapping pairs of point clouds prepared by [36] for registration evaluation.

Performance Tab. 5 shows the point cloud registration performances. Our USIP detector + “Our Desc.” outperforms previous methods with the lowest registration failure rate (Fail %), Relative Translational Error (RTE), Relative Rotation Error (RRE), and highest inlier ratio (Inlier %). In particular, our registration failure rate and inlier ratio are respectively 50% and 2x of the second best keypoint detector + descriptor. We further analyze the performance over different number of RANSAC iterations. The registration failure rate versus the maximum number of RANSAC iterations is shown in Fig. 12. Due to high repeatability, our USIP detector (red line) shows very little drop in performance with decreasing number of RANSAC iterations, while all other algorithms show rapid drops in performances. Additionally, we replace our USIP detector + “Our Desc.” with Random Sampling + “Our Desc.” to demonstrate the effectiveness of our USIP detector. It can be seen from Fig. 12 that the performance of Random Sampling + “Our Desc.” (black line) drops as quickly as other methods with decreasing number of RANSAC iterations.

Effect of USIP Keypoint Saliency Uncertainty Σ on Descriptor Training We show that the keypoint saliency

Method	Oxford					KITTI				
	RTE (m)	RRE ($^{\circ}$)	Fail %	Inlier %	# Iter	RTE (m)	RRE ($^{\circ}$)	Fail %	Inlier %	# Iter
ISS[39] + FPFH[27]	0.40 ± 0.29	1.60 ± 1.02	7.68	8.6	7171	0.33 ± 0.27	1.04 ± 0.77	39.00	8.8	8000
ISS[39] + SI[14]	0.42 ± 0.31	1.61 ± 1.12	12.55	4.7	9888	0.35 ± 0.31	1.11 ± 0.93	41.86	4.6	9401
ISS[39] + USC[30]	0.32 ± 0.27	1.22 ± 0.95	5.98	8.6	7084	0.27 ± 0.28	0.83 ± 0.76	18.62	7.7	8149
ISS[39] + CGF[16]	0.43 ± 0.32	1.62 ± 1.10	12.64	4.9	9628	0.23 ± 0.25	0.69 ± 0.60	8.90	8.4	7670
ISS[39] + 3DMatch[38]	0.49 ± 0.37	1.78 ± 1.21	30.94	5.4	9131	0.30 ± 0.28	0.80 ± 0.67	7.14	8.4	7165
3DFeat-Net[36]	0.30 ± 0.26	1.07 ± 0.85	1.90	13.7	2940	0.26 ± 0.26	0.56 ± 0.46	0.57	12.9	3768
USIP + Our Desc.	0.28 ± 0.26	0.81 ± 0.74	0.93	28.1	523	0.21 ± 0.24	0.42 ± 0.32	0.24	28.0	600

Table 5. Geometric registration performance on Oxford RobotCar and KITTI. The combination of our USIP keypoint detector and “Our Desc.” outperforms existing methods in all criteria with around $2\times$ inlier ratio.

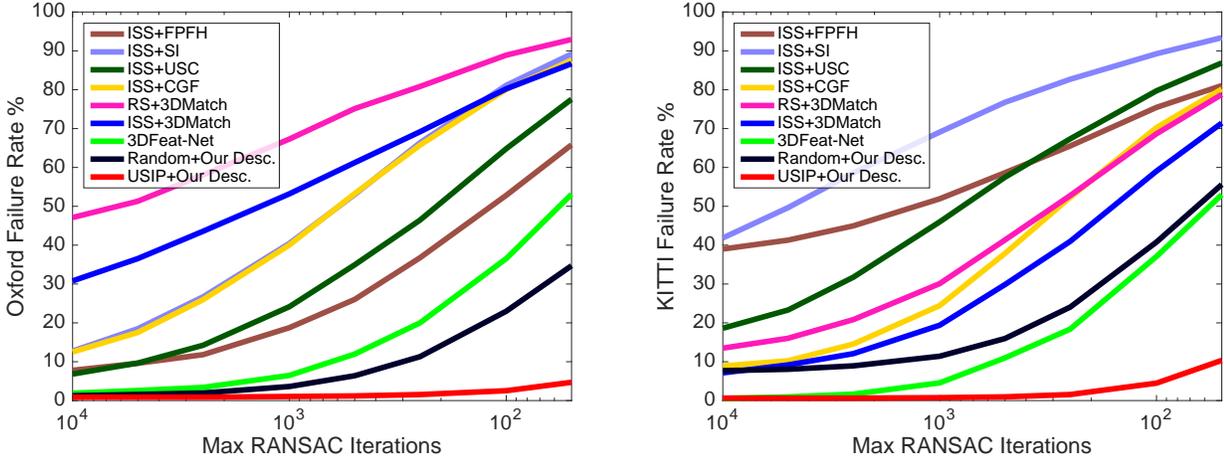


Figure 12. Registration failure rate versus maximum RANSAC iterations in Oxford RobotCar (left) and KITTI (right). Note that the x axis is in logarithmic scale. Our USIP detector + “Our Desc.” (red line) shows very little drop in performance with decreasing number of RANSAC iterations.

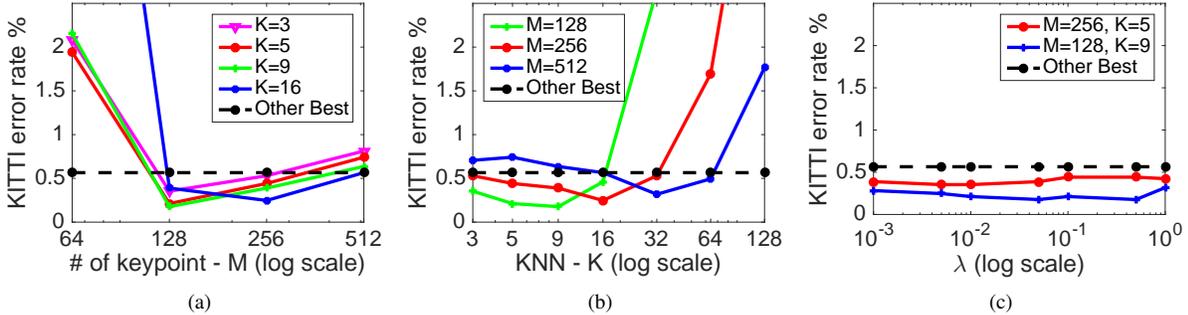


Figure 13. Point cloud registration error rate (%) on KITTI (trained on Oxford). Dash line is the best performance of existing methods. $\lambda = 0.5$ in (a) (b).

uncertainty Σ from our USIP detector improves the performance of “Our Desc.”. To this end, we compare the performances of “Our Desc.” trained with USIP and randomly sampled keypoints, respectively. In particular, the weight w_m from Eq. 15 or Eq. 17 is set to 1 for the randomly sampled keypoints. We denote the descriptor trained with randomly sampled keypoints as “Desc. w. RS”. Tab. 6 shows the registration failure rates of “Desc. w. USIP” and “Desc.

w. RS”. The results show that “Desc. w. USIP” performs better than “Desc. w. RS”, which means that keypoints and saliency uncertainty Σ from our USIP detector improve descriptor training.

Failure %	Oxford		KITTI	
	Desc w. USIP	Desc w. RS	Desc w. USIP	Desc w. RS
USIP	0.93	1.20	0.24	1.02

Table 6. Registration failure rate for “Our Desc.” trained keypoints from our USIP detector and randomly sampled keypoints.

Effect of Parameters M, K, λ We demonstrate the point cloud registration failure rate (%) in Fig. 13, when various USIP detector parameters, M, K, λ , are selected. In Fig. 13 we use the same descriptor mentioned in Sec. D. As shown in Fig. 13, our method outperforms existing methods over a wide range of M, K, λ . We notice our network performance decreases significantly when M is too small or K is too large, *i.e.*, the receptive is too large. This further verifies our design of limiting the receptive field. In addition, Fig. 13 shows that the registration failure rate remains satisfying when λ is small. This is consistent with Fig. 10 that our USIP is able to detect repeatable keypoints even without the point-to-point loss. Nonetheless, it is still important to include the point-to-point loss to ensure that the keypoints are close to the input point cloud.

F. Qualitative Visualization of USIP Keypoints

We show more visualizations of the keypoints detected from our USIP detector on ModelNet40, KITTI, Oxford RobotCar, Redwood in Fig. 16, 17, 18, 19, respectively. NMS and Σ thresholding are applied here. A limitation of our USIP detector is shown in Fig. 16, where there are no or very few keypoints on objects that are highly symmetrical or with smooth surfaces. The saliency uncertainties Σ of the keypoints detected on these objects are large, thus discarded by the Σ thresholding.

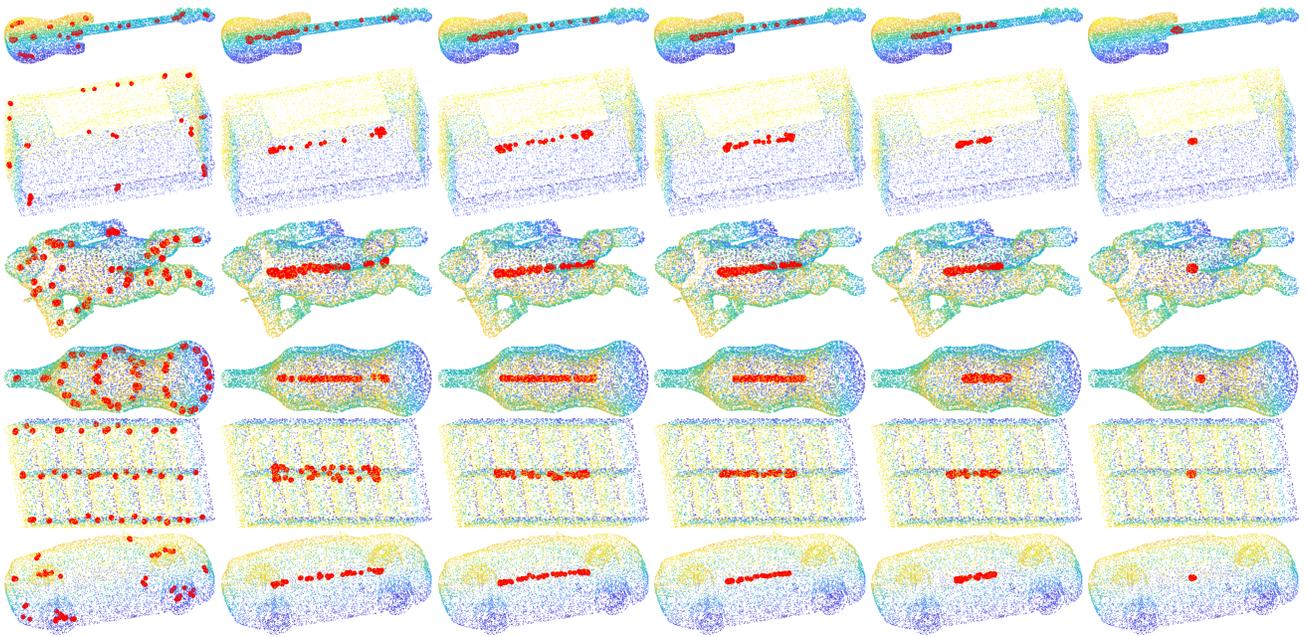


Figure 14. Visualization of FPN degeneracy. $M = 64$ and from left to right: $K = 9, 24, 32, 40, 48, 64$, *i.e.*, receptive field of FPN increases from left to right.

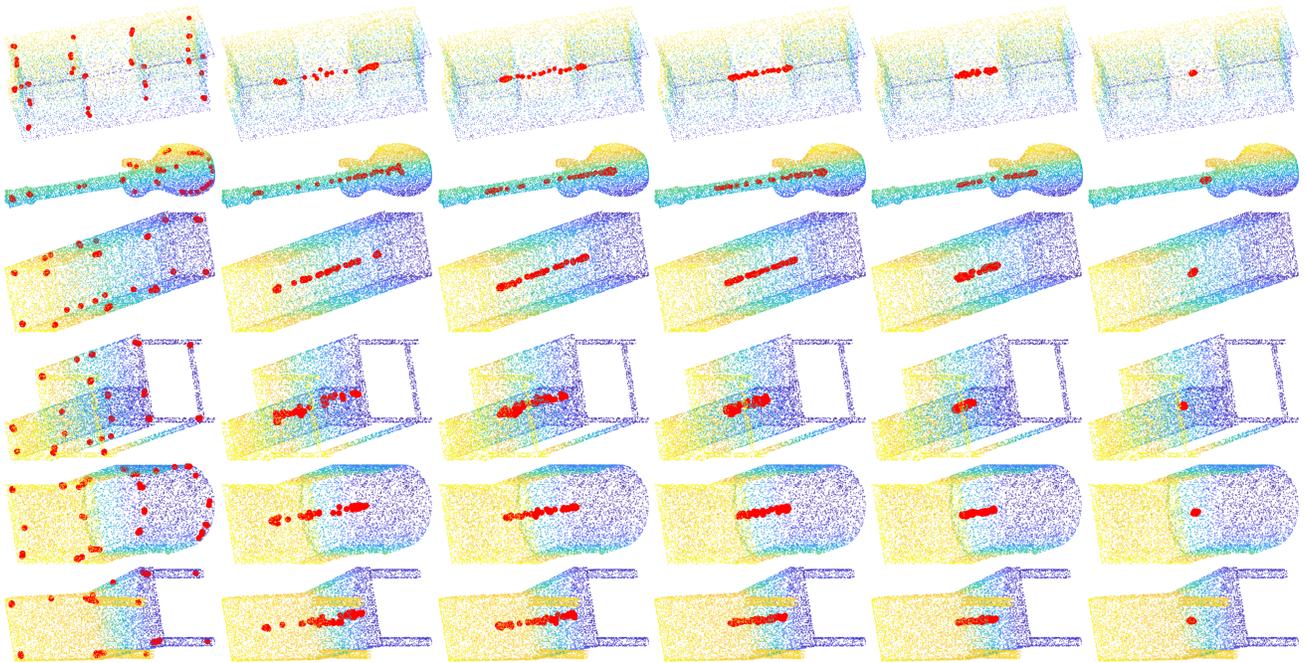


Figure 15. Visualization of FPN degeneracy. $K = 9$ and from left to right: $M = 64, 24, 20, 16, 12, 9$, *i.e.*, receptive field of FPN increases from left to right.

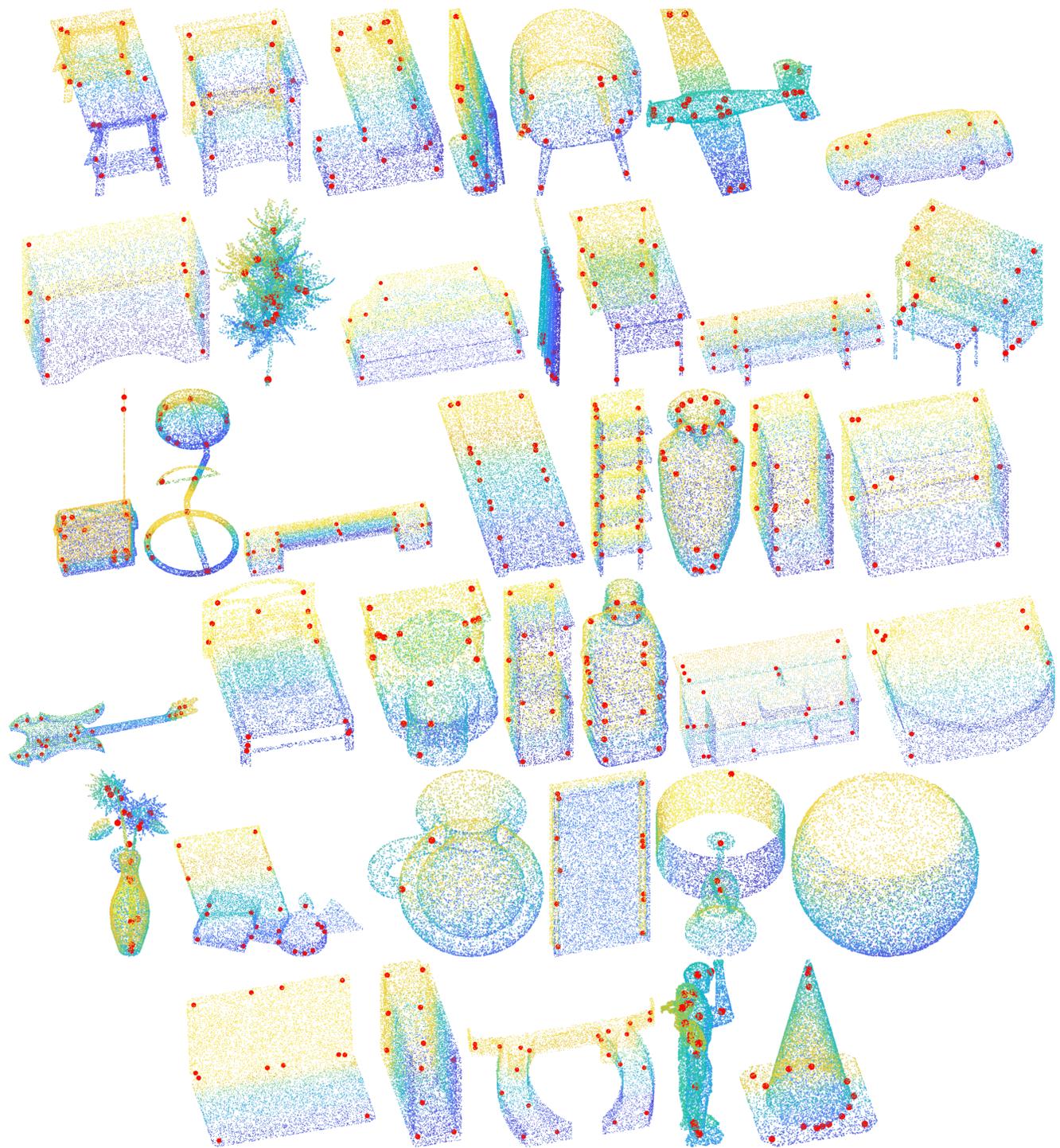


Figure 16. Visualization of USIP keypoints on ModelNet40. Best view with color and zoom-in.

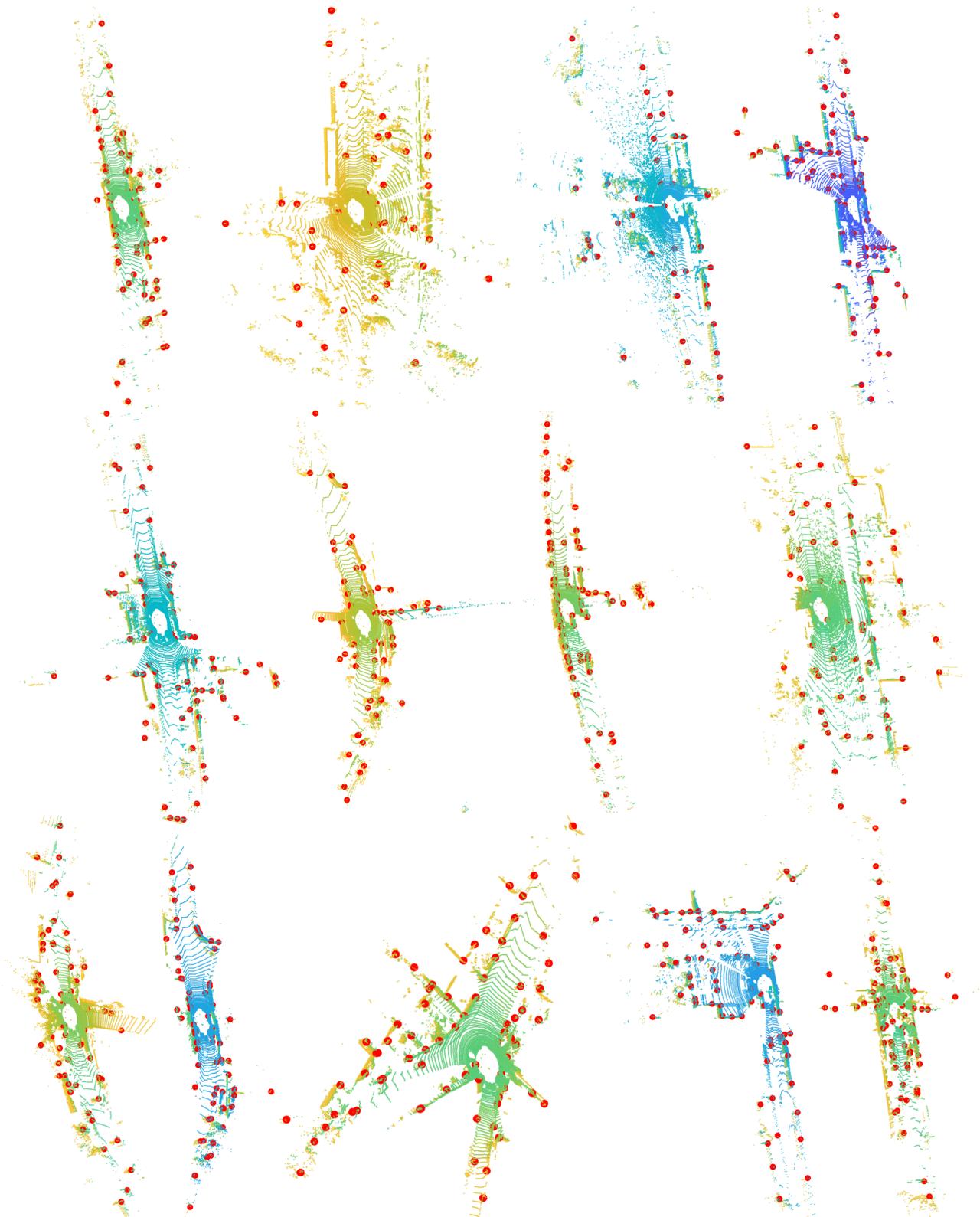


Figure 17. Visualization of USIP keypoints on KITTI with our USIP detector trained on Oxford RobotCar dataset. Best view with color and zoom-in.

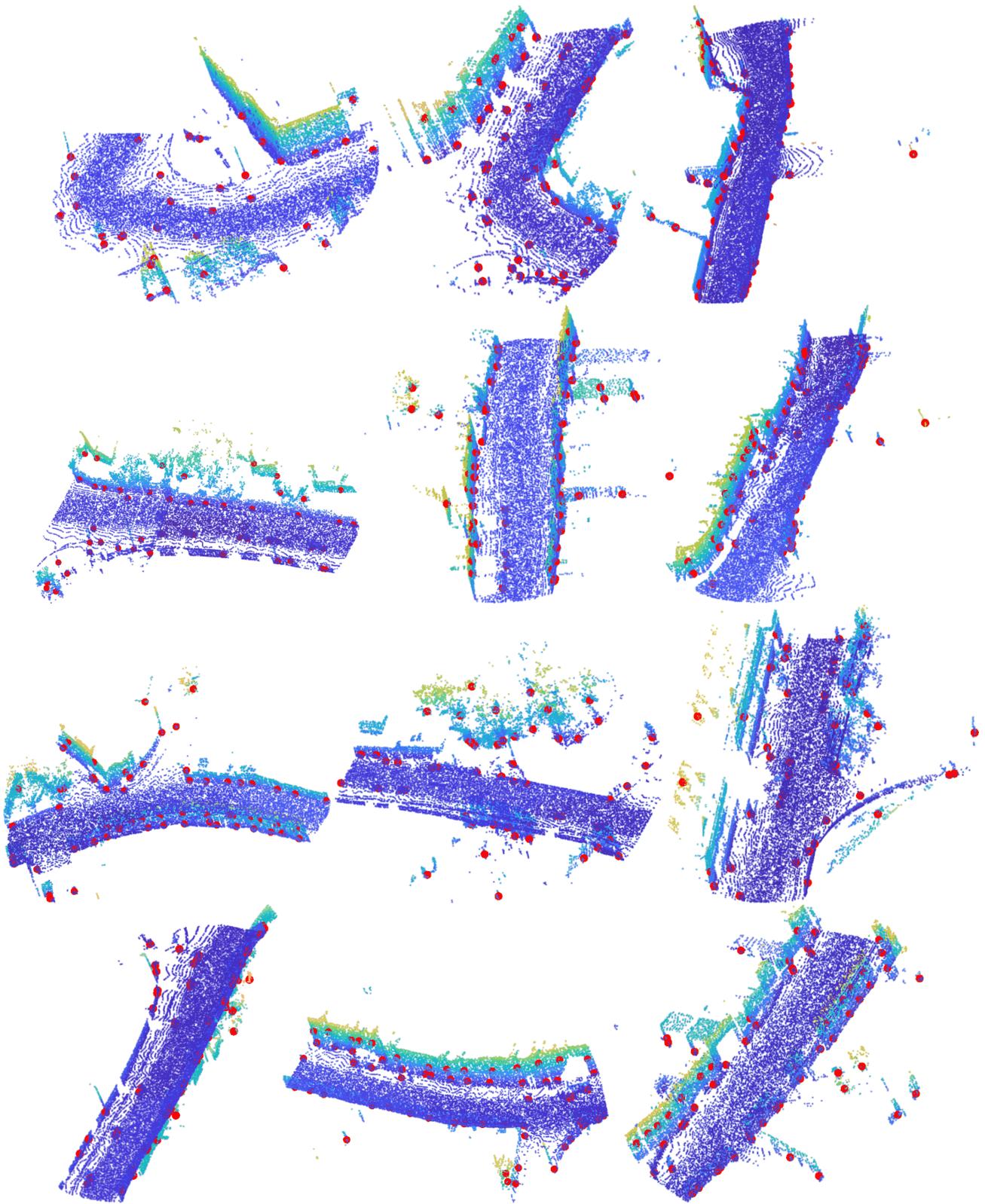


Figure 18. Visualization of USIP keypoints on Oxford RobotCar. Best view with color and zoom-in.

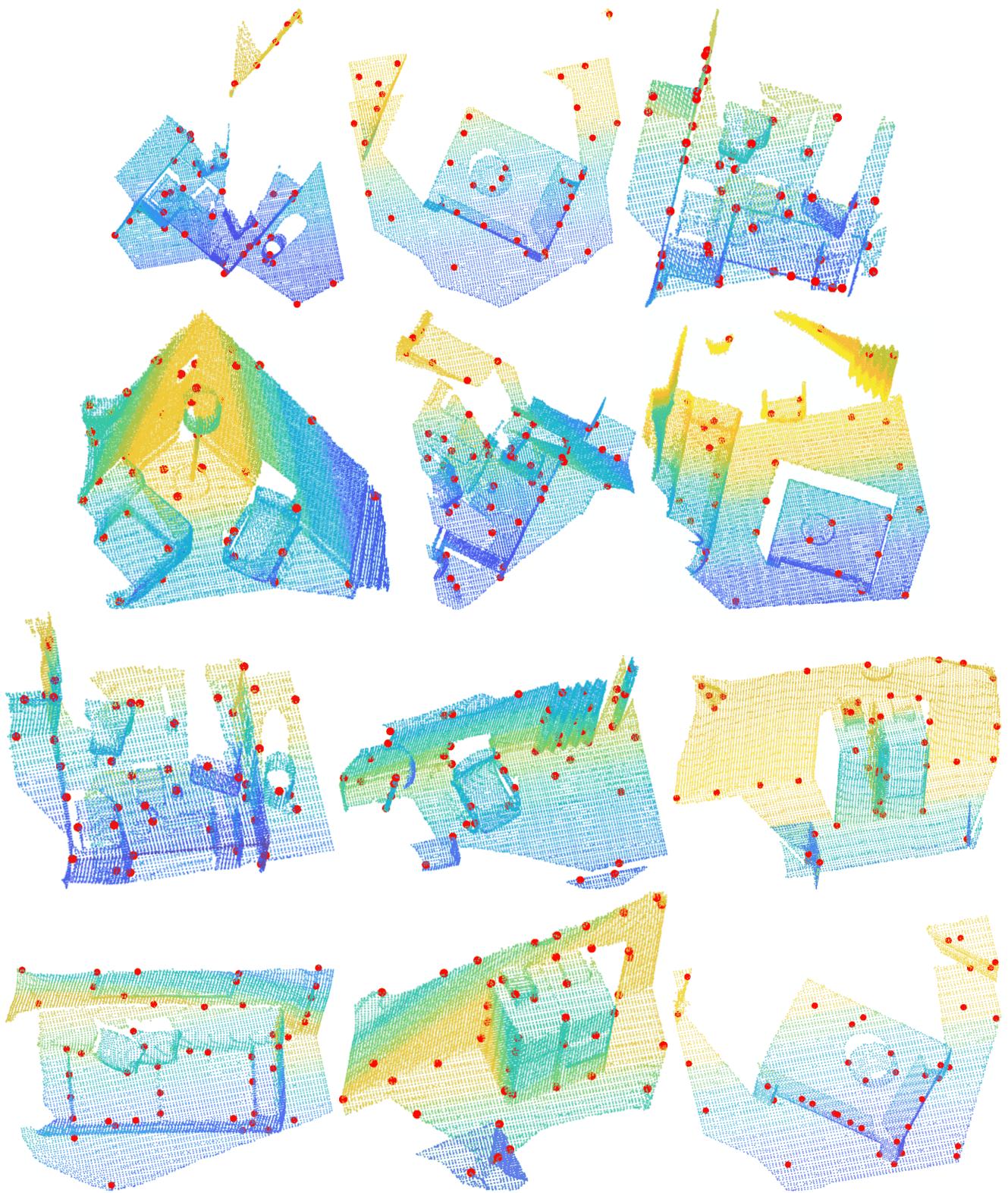


Figure 19. Visualization of USIP keypoints on Redwood with our USIP detector trained on “3D Reconstruction Dataset” [38]. Best view with color and zoom-in.