

High-dimensional Convolutional Networks for Geometric Pattern Recognition

Christopher Choy
NVIDIA

Junha Lee
POSTECH

René Ranftl
Intel Labs

Jaesik Park
POSTECH

Vladlen Koltun
Intel Labs

Abstract

Many problems in science and engineering can be formulated in terms of geometric patterns in high-dimensional spaces. We present high-dimensional convolutional networks (ConvNets) for pattern recognition problems that arise in the context of geometric registration. We first study the effectiveness of convolutional networks in detecting linear subspaces in high-dimensional spaces with up to 32 dimensions: much higher dimensionality than prior applications of ConvNets. We then apply high-dimensional ConvNets to 3D registration under rigid motions and image correspondence estimation. Experiments indicate that our high-dimensional ConvNets outperform prior approaches that relied on deep networks based on global pooling operators.

1. Introduction

Finding structure in noisy data is a general problem that arises in many different disciplines. For example, robust linear regression requires finding a pattern (line, plane) in noisy data. 3D registration of point clouds requires the identification of veridical correspondences in the presence of spurious ones [6]. Structure from motion (SfM) pipelines use verification based on prescribed geometric models to filter spurious image matches [40]. A variety of such applications can benefit from improved methods for the detection of geometric structures in noisy data.

Such detection is challenging. Data points that belong to the sought-after structure often constitute only a small fraction, while the majority are outliers. Various algorithms have been proposed over the years to cope with noisy data [3, 11, 16, 21, 22, 34, 38, 41], but they are usually specific to a subset of problems.

Recent works have advocated for using deep networks [35, 48, 51] to learn robust models to classify geometric structures in the presence of outliers. Deep networks offer significant flexibility and the promise to replace hand-crafted algorithms and heuristics by models learned directly from data. However, due to the unstructured nature of the data in geometric problems, existing works have treated such data as unordered sets, and relied on network

architectures based predominantly on global pooling operators and multi-layer perceptrons (MLPs) [33, 49]. Such network architectures lack the capacity to model local geometric structures and do not leverage the nature of the data, which is often embedded in a (high-dimensional) metric space and has a meaningful geometric structure.

In this work, we introduce a novel type of deep convolutional network that can operate in high dimensions. Our network takes a sparse tensor as input and employs high-dimensional convolutions as the fundamental operator. The distinguishing characteristic of our approach is that it is able to effectively leverage local neighborhood relations together with global context even for high-dimensional data. Our network is fully-convolutional, translation-invariant, and incorporates best practices from the development of ConvNets for two-dimensional image analysis [18, 23, 36].

To demonstrate the effectiveness and generality of our approach, we tackle various geometric pattern recognition problems. We begin with the diagnostic setting of linear subspace detection and show that our construction is effective in high-dimensional spaces and at low signal-to-noise ratios. We then apply the presented construction to geometric pattern recognition problems that arise in computer vision, including registration of three-dimensional point sets under rigid motion (a 6D problem) and correspondence estimation between images under epipolar constraints (a 4D problem). In both settings, the problem is made difficult by the presence of outliers.

Our experiments indicate that the presented construction can reliably detect geometric patterns in high-dimensional data that is heavily contaminated by noise. It can operate in regimes where existing algorithms break down. Our approach significantly improves 3D registration performance when combined with standard approaches for 3D point cloud alignment [6, 39, 53, 54]. The presented high-dimensional convolutional network also outperforms state-of-the-art methods for correspondence estimation between images [48, 51]. All networks and training scripts are available at <https://github.com/chrischoy/HighDimConvNets>.

2. Related Work

Robust model fitting. Fitting a geometric model to a set of observations that are contaminated by outliers is a fundamental problem that frequently arises in computer vision and related fields. The most widely used approach for robust geometric model fitting is *RANdom SAMple Consensus* (RANSAC) [15]. Due to its fundamental importance, many variants and improvements of RANSAC have been proposed over the years [3, 11, 25, 34, 38, 41, 44, 43].

Alternatively, algorithms for robust geometric model fitting are frequently derived using techniques from robust statistics [16, 21, 22, 52], where outlier rejection is performed by equipping an estimator with a cost function that is insensitive to gross outliers. While the resulting algorithms are computationally efficient, they require careful initialization and optimization procedures to avoid poor local optima [53].

Another line of work proposes to find globally optimal solutions to the consensus maximization problem [5, 26, 47]. However, these approaches are currently computationally too demanding for many practical applications.

3D registration. Finding reliable correspondences between a pair of surfaces is an essential step for 3D reconstruction [6, 12, 14, 31]. The problem has been conventionally framed as an energy minimization problem which can be solved using various techniques such as branch and bound [46], Riemannian optimization [37], mixed-integer programming [24], robust error minimization [53], semi-definite programming [20, 28], or random sampling [6].

Recent work has begun to leverage deep networks for geometric registration [2, 30, 32]. These works are predominantly based on PointNet and related architectures, which detect patterns via global pooling operations [33, 49]. In contrast, we develop a high-dimensional convolutional network that operates at multiple scales and can leverage not just global but also local geometric structure.

Image correspondences. Yi *et al.* [48] and Zhang *et al.* [51] reduce essential matrix estimation to classification of correspondences into inliers and outliers. Ranftl and Koltun [35] present a similar formulation for fundamental matrix estimation. Brachmann and Rother [4] propose to learn a neural network that guides hypothesis sampling in RANSAC for model fitting problems. Dang *et al.* [13] propose a numerically stable loss function for essential matrix estimation.

All of these works employ variants of PointNets to classify inliers in an unordered set of putative correspondences. These architectures, based on pointwise MLPs, lack the capacity to model local geometric structure. In contrast, we develop convolutional networks that directly leverage neighborhood relations in the high-dimensional space of correspondences.

3. High-Dimensional Convolutional Networks

In this section, we introduce the two main building blocks of our high-dimensional convolutional network construction: generalized sparse tensors and generalized convolutions.

3.1. Sparse Tensor and Convolution

A tensor is a multi-dimensional array that represents high-order data. A D -th order tensor \mathcal{T} requires D indices to uniquely access its elements. We denote such indices or coordinates as $\mathbf{x} = [x^1, \dots, x^D]$ and the element at the coordinate as $\mathcal{T}[\mathbf{x}]$ similar to how we access components in a matrix. Likewise, a sparse tensor is a high-dimensional extension of a sparse matrix where the majority of the elements are 0. Concretely,

$$\mathcal{T}[\mathbf{x}_i] = \begin{cases} \mathbf{f}_i & \text{if } \mathbf{x}_i \in \mathcal{C} \\ 0 & \text{otherwise,} \end{cases} \quad (1)$$

where $\mathcal{C} = \{\mathbf{x}_i \mid \mathbf{x}_i \in \mathbb{N}^D, \mathcal{T}[\mathbf{x}_i] \neq \mathbf{0}\}_{i=1}^N$ is the set of coordinates with non-zero values, N is the number of non-zero elements, and \mathbf{f}_i is the non-zero value at the i -th coordinate. A sparse tensor feature map is a $(D + 1)$ -th order tensor with $\mathbf{f}_i \in \mathbb{R}^{N_{D+1}}$ as we use the last dimension to denote the feature dimension. A sparse tensor has the constraint that $\mathbf{x}_i \in \mathbb{N}^D$. We extend the sparse tensor coordinates to integer indices $\mathbf{x}_i \in \mathbb{Z}^D$ and define $\mathcal{T} \in \mathbb{R}^{\mathbb{N}_0^D \times N_{D+1}}$ where \mathbb{N}_0 denotes the cardinality of the integer space $|\mathbb{Z}|$ to define a generalized sparse tensor.

A convolution on this generalized sparse tensor can then be defined as a simple extension of the generalized sparse convolution [7]:

$$\mathbf{f}_{\mathbf{x}}^{\text{out}} = \sum_{\mathbf{i} \in \mathcal{N}^D(\mathbf{x}) \cap \mathcal{C}^{\text{in}}} \mathbf{W}_{\mathbf{i}} \mathbf{f}_{\mathbf{x}+\mathbf{i}}^{\text{in}} \quad \text{for } \mathbf{x} \in \mathcal{C}^{\text{out}}, \quad (2)$$

where \mathcal{C}^{in} and \mathcal{C}^{out} are the set of input and output locations which is predefined by the user, $\mathbf{W}_{\mathbf{i}}$ is a weight matrix, and $\mathcal{N}^D(\mathbf{x})$ defines a set of neighbors of \mathbf{x} which is defined by the shape of the convolution kernel. For example, if the convolution kernel is a hypercube of size K , $\mathcal{N}^D(\mathbf{x}) \cap \mathcal{C}^{\text{in}}$ is a set of all the non-zero elements of the input sparse tensor centered at \mathbf{x} within the L_{∞} -ball of extent K .

3.2. Convolutional Networks

We design a high-dimensional fully-convolutional neural network for sparse tensors (sparse tensor networks) based on generalized convolution [7, 9]. We use U-shaped networks [36] to capture large receptive fields while maintaining the original resolution at the final layer. The network has residual connections [18] within layers with the same resolution and across the network to speed up convergence

and to recover the lost spatial resolution in the last layer. The network architecture is illustrated in Fig. 1.

For computational efficiency in high dimensions, we use the cross-shaped kernel [7] for all convolutions. We denote the kernel size by K . The cross-shaped kernel has non-zero weights only for the $K - 1$ nearest neighbors along each axis, which results in one weight parameter for the center location and $K - 1$ weight parameters for each axis. Note that a cross-shaped kernel is similar to separable convolution, where a full convolution is approximated by D one-dimensional convolutions of size K . Both types of kernels are rank-1 approximations of the full hyper-cubic kernel K^D , but separable convolution requires KD matrix multiplications, whereas the cross-shaped kernel requires only $(K - 1)D + 1$ matrix multiplications.

3.3. Implementation

We extend the implementation of Choy *et al.* [7], which supports arbitrary kernel shapes, to high-dimensional convolutional networks. To implement the sparse tensor networks, we need an efficient data structure that can generate a new sparse tensor as well as find neighbors within the sparse tensor. Choy *et al.* [7] use a hash table that is efficient for both insertion and search. We replaced the hash table with a faster and more efficient variant [1]. In addition, as the neighbor search can be run in parallel, we create an iterator function that can run in parallel with OpenMP [29] by dividing the table into smaller parallelization blocks.

Lastly, U-shaped networks generate hierarchical feature maps that expand the receptive field. Choy *et al.* [7] use stride- K convolutions with kernel size $\geq K$ to generate lower-resolution hierarchical feature maps. Still, such implementation requires iterating over at least K^D elements within a hypercubic kernel as the coordinates are stored in a hash table, which results in $O(NK^D)$ complexity where N is the cardinality of input. Consequently, it becomes infeasible to store the weights on the GPU for high-dimensional spaces. Instead of strided convolutions, we propose an efficient implementation of stride- K sum pooling layers with kernel size K . Instead of iterating over all possible neighbors, we iterate over all input coordinates and round them down to multiples of K , which requires only $O(N)$ complexity.

4. Geometric Pattern Recognition

Our approach uses convolutional networks to recognize geometric patterns in high-dimensional spaces. Specifically, we classify each point \mathbf{x}_i in a high-dimensional dataset $\mathcal{X} = \{\mathbf{x}_i\}_{i=1}^N$ as an inlier or an outlier. We start by validating our approach on synthetic datasets of varying dimensionality and then show results on 3D registration and essential matrix estimation.

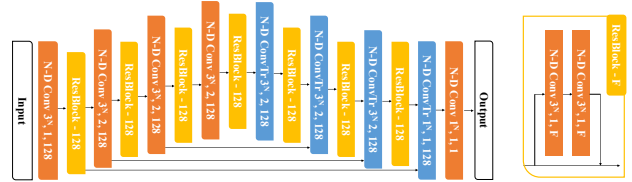


Figure 1: A generic U-shaped high-dimensional convolutional network architecture. The numbers next to each block indicate kernel size, stride, and the number of channels. The strided convolutions that reduce the resolution of activations are shifted upward to indicate different levels of resolution.

For all experiments, we first quantize the input coordinates to create a sparse tensor of order $D + 1$, where the last dimension denotes the feature channels. The network then predicts a logit score for each non-zero element in the sparse tensor to indicate if a point is part of the geometric pattern or if it is an outlier.

4.1. Line and Plane Detection

We first test the capabilities of our fully-convolutional networks on simple high-dimensional pattern recognition problems that involve detecting linear subspaces amidst noise. Our dataset consists of uniformly sampled noise from the D -dimensional space and a small number of samples from a line under a Gaussian noise model. The number of outliers increases exponentially in the dimension ($O(L^D)$), while the number of inliers increases sublinearly ($O(\sqrt{LD})$), where L is the extent of the domain. Further details are given in the supplement. The network predicts a likelihood score for each non-zero element in the input sparse tensor, and we threshold inliers with probability ≥ 0.5 . We estimate the line equation from the predicted inliers using unweighted least squares.

We use PointNet variants as the baselines for this experiment [33, 48, 49]. For Zaheer *et al.* [49], we were not able to get reasonable results with the network architecture proposed in the paper. We thus augmented the architecture with batch normalization and instance normalization layers after each linear transformation similar to Yi *et al.* [48], which boosted performance significantly. For all experiments, we use the cross-entropy loss. We used the same training hyperparameters, including loss, batch size, optimizer, and learning rate schedule for all approaches.

We use three metrics to analyze the performance of the networks: Mean Squared Error (MSE), F1 score, and Average Precision (AP). For the MSE, we estimate the line equation with Least Squares to fit the line to the inliers. The second metric is the F1 score, the harmonic mean of precision and recall. In many problems, F1 score is a direct indicator of the performance of a classifier, and we also found a

Table 1: Line detection in high-dimensional spaces in the presence of extreme noise. All networks are trained with the cross-entropy loss for 40 epochs. Inlier ratios are listed on the left. In the 32-dimensional setting, only 7 out of 10,000 data points are inliers. The table reports Mean Squared Error (MSE), F1 score, and Average Precision (AP) of our approach (a high-dimensional ConvNet) versus baselines (PointNet variants). MSE: lower is better. F1 and AP: higher is better.

Dim.	Inlier Ratio	Qi <i>et al.</i> [33]			Zaheer <i>et al.</i> [49] + BN + IN			Yi <i>et al.</i> [48]			Ours		
		MSE	F1	AP (AUC)	MSE	F1	AP (AUC)	MSE	F1	AP (AUC)	MSE	F1	AP (AUC)
4	15.59%	1.337	0.025	0.164	6.33E-4	0.867	0.936	9.11E-5	0.981	0.996	2.33E-5	0.998	0.999
8	5.54%	2.369	0.022	0.065	0.001	0.891	0.955	2.45E-4	0.946	0.989	1.64E-5	0.999	0.999
16	2.75%	3.854	0.012	0.034	4.86E-4	0.970	0.992	0.002	0.962	0.986	3.39E-5	0.999	0.999
24	0.40%	5.372	0.021	0.011	0.676	0.634	0.691	0.775	0.610	0.674	5.34E-5	0.994	0.996
32	0.07%	6.715	0.012	6.71E-5	-	0.0	0.295	-	0.0	0.050	0.010	0.669	0.689

Table 2: Plane detection in high-dimensional spaces in the presence of extreme noise. Inlier ratios are listed on the left. In the 32-dimensional setting, fewer than 5 out of 100,000 data points are inliers. The table reports the F1 score and Average Precision (AP) of our approach (a high-dimensional ConvNet) versus baselines (PointNet variants). Higher is better.

Dim.	Inlier Ratio	Qi <i>et al.</i> [33]		Zaheer <i>et al.</i> [49] + BN + IN		Yi <i>et al.</i> [48]		Ours	
		F1	AP (AUC)	F1	AP (AUC)	F1	AP (AUC)	F1	AP (AUC)
4	29.96%	0.0	0.315	0.980	0.996	0.993	0.999	0.991	0.998
8	8.07%	0.0	0.088	0.985	0.999	0.990	0.999	0.998	0.999
16	0.34%	0.0	0.004	0.155	0.299	0.182	0.359	0.951	0.961
24	0.01%	0.0	1.61E-4	0.032	0.133	0.0	0.081	0.304	0.346
32	4.64E-3%	0.0	5.56E-5	0.0	0.221	0.0	0.023	0.138	0.240

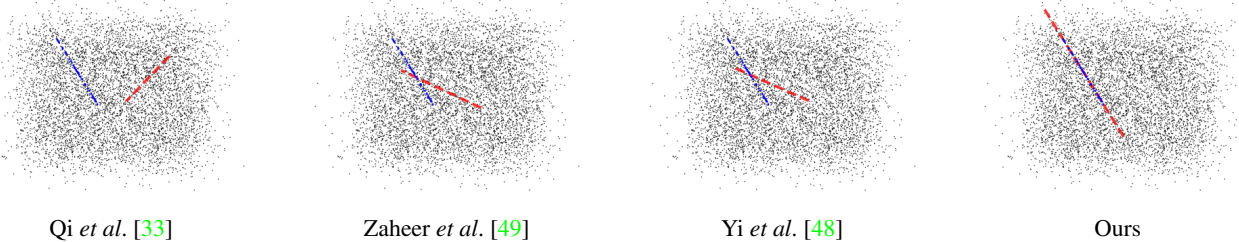


Figure 2: 16D line detection projected to a 2D plane for visualization. Black dots are noise and blue dots are samples from a line in the 16D space. The dashed red line is the prediction of the respective method. Samples from the ground-truth line (blue) are enlarged by a factor of 10 for visualization.

strong correlation between F1 score and the mean squared error. The final metric we use is average precision (AP), which measures the area under the precision-recall curve. We report the results in Tab. 1 and provide qualitative examples in Fig. 2. Tab. 1 lists the inlier ratio to indicate the difficulty of each task.

In a second experiment, we create another synthetic dataset where the inlier pattern is sampled from a plane spanned by two vectors: $\{c_1 \mathbf{v}_1 + c_2 \mathbf{v}_2 + \mathbf{c} \mid c_1, c_2 \in \mathbb{R}\}$. The two basis vectors are sampled uniformly from the D -dimensional unit hypercube. We use the same training procedure for our network and the baselines, and report the results in Tab. 2.

We found that the convolutional network is more robust to noise in high-dimensional spaces than the PointNet variants. In addition, the convolutional network training converges quickly, as shown in Fig. 3, which is a further indication that the architecture can effectively leverage the struc-

ture of the data.

4.2. 3D Registration

A typical 3D registration pipeline consists of 1) feature extraction, 2) feature matching, 3) match filtering, and 4) global registration. In this section, we show that in the *match filtering* stage the correct (inlier) correspondences form a 6-dimensional geometric structure. We then extend our geometric pattern recognition networks to identify inlier correspondences in this 6-dimensional space.

Let \mathcal{X} be a set of points sampled from a 3D surface, $\mathcal{X} = \{\mathbf{x}_i \mid \mathbf{x}_i \in \mathbb{R}^3\}_{i=1}^N$, and let \mathcal{X}' be a subset of \mathcal{X} that went through a rigid transformation T , $\mathcal{X}' = \{T(\mathbf{x}) \mid \mathbf{x} \in \mathcal{S}, \mathcal{S} \subseteq \mathcal{X}\}$. For example, \mathcal{X}' could be a 3D scan from a different perspective that has an overlap with \mathcal{X} . We denote a correspondence between points $\mathbf{x}_i \in \mathcal{X}$ and $\mathbf{x}'_j \in \mathcal{X}'$ as $\mathbf{x}_i \leftrightarrow \mathbf{x}'_j$. When we form an ordered pair $(\mathbf{x}_i, \mathbf{x}'_j) \in \mathbb{R}^6$, the ground truth correspondences satisfy $T(\mathbf{x}) = \mathbf{x}'$ along the

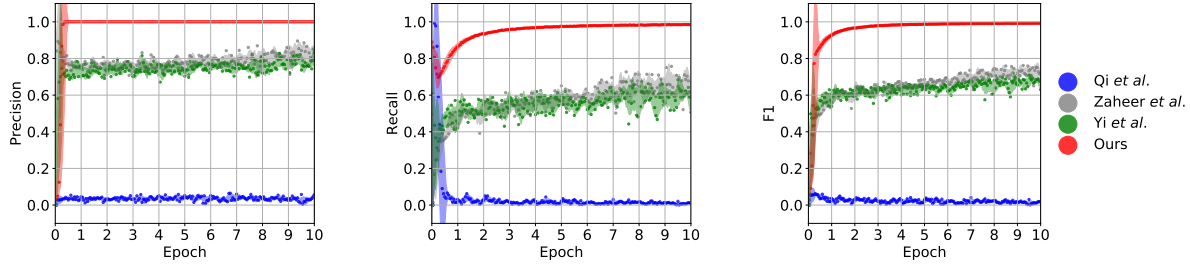


Figure 3: 16D line detection, progression of training. We plot the running mean and standard deviation of precision, recall, and F1 score on the validation set. Our high-dimensional convolutional network quickly attains much higher accuracy than the baselines.

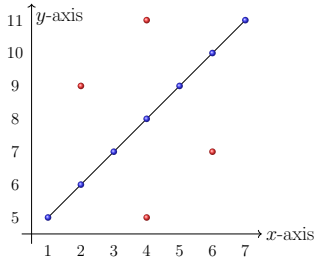


Figure 4: The set \mathcal{Y} is rigid translation of the set $\mathcal{X} = [1, \dots, 7]$, $\mathcal{Y} = \{x+4|x \in \mathcal{X}\}$. The ordered pairs of correspondences (blue) form a line segment while outliers (red) form random noise outside the line.

common 3D geometry \mathcal{S} whereas an incorrect correspondence implies $T(\mathbf{x}) \neq \mathbf{x}'$. For example, in Fig. 4, we visualize the ordered pairs from 1-dimensional sets. Note that the inliers follow the geometry of the inputs and form a line segment. Similarly, the geometry $(\mathbf{x}, \mathbf{x}') \in \mathbb{R}^6$ or $(\mathbf{x}, T(\mathbf{x}))$ for $\mathbf{x} \in \mathcal{S}$ forms a surface in 6-dimensional space.

Theorem 1 *The 6D representation of ground-truth 3D correspondences $(\mathbf{x}, \mathbf{x}')$ lies on the intersection of three hyperplanes where each hyperplane is a row of the following block equation $\begin{bmatrix} R & -I \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \mathbf{x}' \end{bmatrix} + \mathbf{t} = \mathbf{0}$.*

Proof: A ground-truth 3D correspondence, $\mathbf{x} \leftrightarrow \mathbf{x}'$, must satisfy $R\mathbf{x} + \mathbf{t} = \mathbf{x}'$. We move \mathbf{x}' to the LHS and convert both \mathbf{x} and \mathbf{x}' to the 6D representation to get three hyperplane equations. \square

Corollary 1.1 *The rank of the block matrix, $\begin{bmatrix} R & -I \end{bmatrix}$, is 3 since R and I are orthogonal matrices. Thus, all hyperplanes defined by the block matrix intersect each other and the 6D inlier correspondences form a 3D plane since the solution of the intersection of three hyperplanes in the 6D space, $\begin{bmatrix} R & -I \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \mathbf{x}' \end{bmatrix} + \mathbf{t} = \mathbf{0}$, is a 3D plane.*

We can thus use our high-dimensional convolutional network construction to segment the 6-dimensional set of cor-

respondences into inliers and outliers by estimating the inlier likelihood for each correspondence.

Network. We use a 6-dimensional instantiation of the U-shaped convolutional network presented in Sec. 3. As the dimensionality is manageable, we use hypercubic kernels. The network takes an order-6 sparse tensor whose coordinates are correspondences $(\mathbf{x}_i, \mathbf{x}'_j) \in \mathbb{R}^6$. We discretize the coordinates with the voxel size used to extract features. Our baseline is Yi *et al.* [48], which takes dimensionless mean-centered correspondences without discretization. We train the networks to predict the inlier probability of each correspondence with the balanced cross-entropy loss.

Dataset. We use the 3DMatch dataset for this experiment [50]. The 3DMatch dataset is a composition of various 3D scan datasets [17, 45, 50] and thus covers a wide range of scenes and different types of 3D cameras. We integrate RGB-D images to form fragments of the scenes following [50]. During training, we randomly rotate each scene on the fly to augment the dataset. We use a popular hand-designed feature descriptor, FPFH [39], to compute correspondences. Note, however, that our pipeline is agnostic to the choice of feature and can also be used with learned features [9].

We follow the standard procedures in the 3D registration literature to generate candidate correspondences. First, since 3D scans often exhibit irregular densities, we resample the input point clouds using a voxel grid to produce a regular point cloud. We use voxel sizes of 2.5cm and 5cm for our experiments. Next, we compute FPFH features and find the nearest neighbor for each point in feature space to form correspondences. The correspondences obtained from this procedure often exhibit a very low inlier ratio, as little as 0.87% with a 2.5cm voxel size. Among these correspondences, we regard $\mathbf{x} \leftrightarrow \mathbf{x}'$ as an inlier if it satisfies $\|\mathbf{T}(\mathbf{x}) - \mathbf{x}'\|_2 < \tau$ and all others as outliers. We set τ to be two times the voxel size.

Finally, we use a registration method to convert the filtered correspondences into a final registration result. We show results with two different registration methods. The

Table 3: Pairwise registration on 3DMatch test scenes with 2.5cm downsampling. Translation Error (TE), Rotation Error (RE), success rate. Registration is considered successful if TE < 30cm and RE < 15°.

	Inlier Ratio	FPFH + FGR			FPFH + Ours + FGR			FPFH + RANSAC			FPFH + Ours + RANSAC		
		TE	RE	Succ. Rate	TE	RE	Succ. Rate	TE	RE	Succ. Rate	TE	RE	Succ. Rate
Kitchen	1.62%	10.98	4.99	37.15	5.68	2.21	65.61	6.25	2.17	44.47	5.90	1.98	69.57
Home 1	2.71%	11.12	4.40	45.51	6.52	2.08	80.77	7.07	2.19	61.54	6.00	1.87	80.13
Home 2	2.83%	9.61	3.83	36.54	7.13	2.56	64.42	6.47	2.40	50.00	7.86	2.56	69.71
Hotel 1	1.35%	12.31	5.09	33.19	7.95	2.65	76.11	7.48	2.75	48.67	7.38	2.38	80.09
Hotel 2	1.54%	12.27	5.22	25.00	7.86	2.56	69.23	9.54	3.18	47.12	6.40	2.25	70.19
Hotel 3	1.59%	13.52	7.04	27.78	5.39	1.99	72.22	5.91	2.46	59.26	5.85	2.36	81.48
Study	0.87%	16.10	6.01	16.78	9.61	2.64	53.42	10.05	3.01	30.48	8.51	2.23	56.16
Lab	1.59%	10.48	4.80	42.86	7.69	2.44	61.04	8.01	2.31	45.45	6.64	2.12	68.83
Average		12.05	5.17	33.10	7.23	2.39	67.85	7.60	2.56	48.37	6.82	2.22	72.02

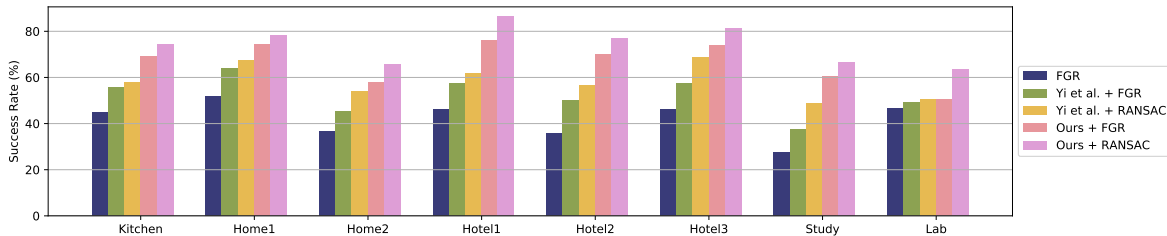


Figure 5: The success rate of baseline methods and ours on the 3DMatch benchmark [50] with 5cm voxel size. FGR denotes registration with FPFH [39] and Zhou *et al.* [53], Yi *et al.* + X denotes FPFH filtering with Yi *et al.* [48] and registration with X, and Ours + X denotes our method for filtering followed by registration with X.

first is Fast Global Registration [53], which directly minimizes a robust error metric. The second is a variant of RANSAC [15] that is specialized to 3D registration [54].

Evaluation. We use three standard metrics to evaluate registration performance: rotation error, translation error, and success rate. The rotation error measures the absolute angular deviation from the ground truth rotation $\hat{\mathbf{R}}$, $\arccos \frac{\text{Tr}(\hat{\mathbf{R}}^T \mathbf{R}) - 1}{2}$. Similarly, the translation error measures the deviation of the translation $\|\hat{\mathbf{t}} - \mathbf{t}\|_2$. When we report these metrics, we exclude alignments that exceed a threshold following [9] since the results of the registration methods [53, 15] can be arbitrarily bad when registration fails. Finally, the success rate is the ratio of registrations that were successful; registration is considered successful if both rotation and translation errors are within the respective thresholds. For all experiments, we use a rotation error of 15 degrees and a translation error of 30cm as the thresholds.

Tab. 3 shows the 3D registration pipelines with and without our network to filter the outliers. Note that for FGR [53], we observe a considerable improvement with our network since FGR assumes more accurate correspondences as inputs. The improvement is smaller with RANSAC, since it is more robust to high outlier rates. Although the inlier ratio is as low as 1% for many 3D scene pairs, our network generates very accurate predictions. Similar to the linear

regression experiments in Fig. 3, we find that the network converges very quickly. We compare to the model of Yi *et al.* [48] in Fig. 5 with 5cm voxel size to study the robustness of the 6-dimensional convolutional network to voxel size and find that the convolutional network improves the registration success rate significantly even for the higher inlier ratio seen with a 5cm discretization resolution. Qualitatively, our network accurately filters out outliers even in the presence of extreme noise (Fig. 6).

4.3. Filtering Image Correspondences

In this section, we apply the high-dimensional convolutional network to image correspondence inlier detection. In the projective space \mathbb{P}^2 , an inlier correspondence $\mathbf{u} \leftrightarrow \mathbf{u}'$ must satisfy $\mathbf{u}'^T \mathbf{E} \mathbf{u} = 0$, where \mathbf{E} is an essential matrix, \mathbf{u} denotes a normalized homogeneous coordinate $\mathbf{u} = \mathbf{K}^{-1} \mathbf{x}$, \mathbf{x} is the corresponding homogeneous image coordinate, and \mathbf{K} is the camera intrinsic matrix. When we expand $\mathbf{u}'^T \mathbf{E} \mathbf{u} = 0$, we get $u'^1 A u^1 + u'^2 B u^1 + u'^1 C u^2 + u'^2 D u^2 + E u'^1 + F u'^2 + G u^1 + H u^2 + I = 0$, which is a quadri-variate quadratic function. If there is a real-valued solution, there are infinitely many solutions that form either an ellipse (sphere), a parabola, or a hyperbola. These are known as conic sections. Thus a set of ground-truth image correspondences will form a hyper conic section in 4-dimensional space. We use a convolutional network to predict the likeli-

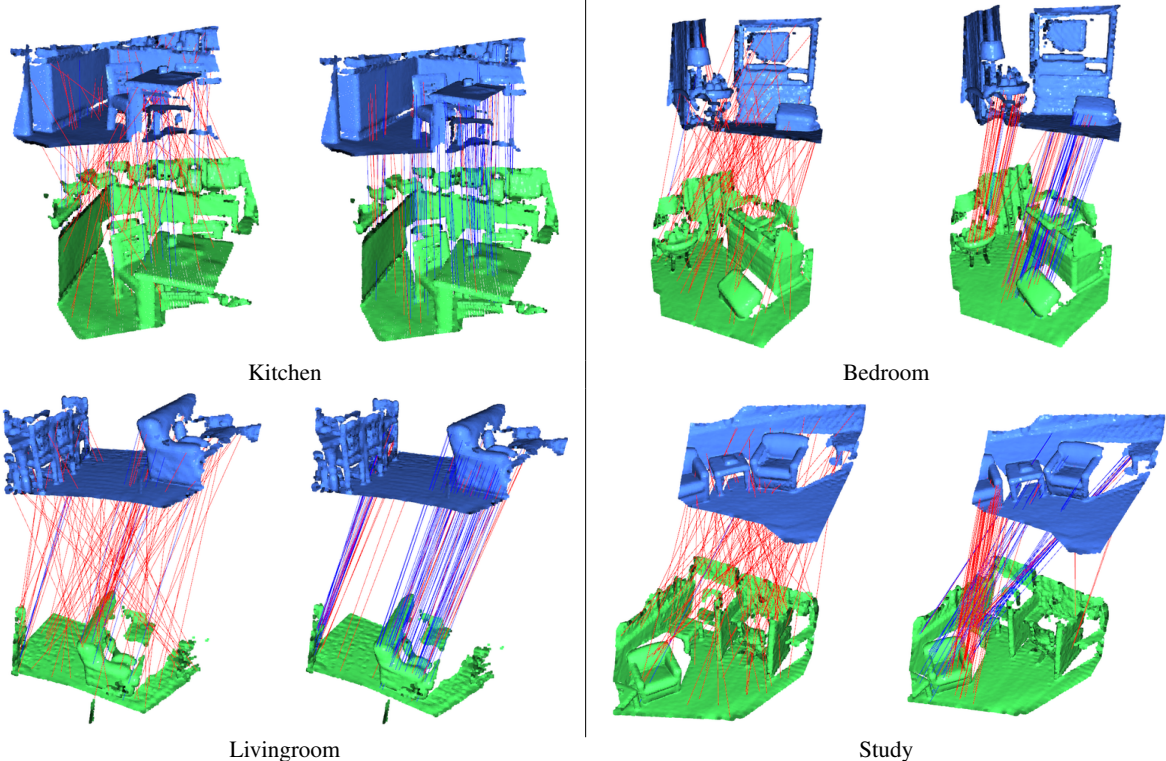


Figure 6: Visualization of color-coded correspondences before and after outlier filtering (Sec. 4.2). For each pair, we visualize 100 random correspondences from the candidate set on the left and 100 random correspondences after outlier pruning on the right. Red lines are outlier correspondences and blue lines are inlier correspondences. On the bottom right pair, there are two identical chairs. The average inlier ratio is 1.76%.

hood that a correspondence is an inlier.

Dataset: YFCC100M. We use the large-scale photo-tourism dataset YFCC100M [42] for the experiment. The dataset contains 100M Flickr images of tourist hot-spots with metadata, which is curated into 72 locations with camera extrinsics estimated using SfM [19]. We follow Zhang *et al.* [51] to generate a dataset and use 68 locations for training and the others for testing. We filtered any image pairs that have fewer than 100 overlapping 3D points from SfM to guarantee non-zero overlap between images.

We use SIFT features [27] to create correspondences and label a correspondence to be a ground-truth inlier if the symmetric epipolar distance of the correspondence is below a certain threshold using the provided camera parameters, *i.e.*,

$$\left(\frac{r^2}{l_1^2 + l_2^2} + \frac{r^2}{l_1'^2 + l_2'^2} \right) < \tau, \quad (3)$$

where $l = \mathbf{u}'^\top \mathbf{E} = (l_1, l_2, l_3)$ is a homogeneous line, $l' = \mathbf{E}\mathbf{u}$, and $r = \mathbf{u}'^\top \mathbf{E}\mathbf{u}$.

Network. We convert a set of candidate image correspondences into an order-5 sparse tensor with four spatial dimensions and vectorized features. The coordinates are defined

as normalized image coordinates \mathbf{u} . We define integer coordinates by discretizing the normalized image coordinates with quantization resolution 0.01 and additionally use the normalized coordinates as features. We use state-of-the-art baselines on the YFCC dataset and two variants of convolutional networks for this task. The first variant is a U-shaped convolutional network (Ours); the second network is a ResNet-like network with spatial correlation modules (Ours + SC). Spatial correlation modules [51] are blocks of shared MLPs that encode global context from a set of correspondences.

Note that unlike the FPFH descriptor [39], which extracts features densely, SIFT features are extracted from only a few keypoints sparsely in the image. In the 4-dimensional space of correspondences, the sparsity gets even worse as the volume increases multiplicatively while the number of points (correspondences) stay the same. Such sparsity leads to fewer neighbors in the high-dimensional space, which results in the degeneration of the convolution to a multi-layer perceptron. Our convolutional networks remain effective in this setting, but their distinctive ability to leverage local geometric structure is not strongly utilized. To increase the density of the correspondences in

the 4-dimensional space, we use a dense fully convolutional feature UCN [10]. We train the UCN on the YFCC100M training set for 100 epochs and we follow the preprocessing stage outlined in the open-source version [8] to create correspondences. Training and testing on UCN follow the same standard procedure. We use the balanced cross-entropy loss [48] for all experiments.

Evaluation. We use precision, recall, and F1 score to evaluate correspondence classification accuracy. We use $\tau = 10^{-4}$ for the distance threshold to define ground truth-correspondences. Our network predicts a correspondence to be accurate if its inlier probability prediction is above 0.5. We provide quantitative results in Tab. 4 and qualitative results in Fig. 7. Our approaches outperform the PointNet variants [48, 51] as measured by the F1 score.

5. Conclusion

Many interesting problems in computer vision involve geometric patterns in high-dimensional spaces. We proposed high-dimensional convolutional networks for geometric pattern recognition. We presented a fully-convolutional network architecture that is efficient and able to find patterns in high-dimensional data even in the presence of severe noise. We validated the efficacy of our approach on tasks such as line and plane detection, 3D registration, and geometric filtering of image correspondences.

Acknowledgements

This work was partially supported by a National Research Foundation of Korea (NRF) grant and funded by the Korean government (MSIT) (No. 2020R1C1C1015260).

References

- [1] Martin Ankerl. Robin hood hashing. <https://github.com/martinus/robin-hood-hashing>, 2019. 3
- [2] Yasuhiro Aoki, Hunter Goforth, Rangaprasad Arun Srivatsan, and Simon Lucey. PointNetLK: Robust & efficient point cloud registration using pointnet. In *CVPR*, 2019. 2
- [3] Eric Brachmann, Alexander Krull, Sebastian Nowozin, Jamie Shotton, Frank Michel, Stefan Gumhold, and Carsten Rother. DSAC - differentiable RANSAC for camera localization. In *CVPR*, 2017. 1, 2
- [4] Eric Brachmann and Carsten Rother. Neural-guided RANSAC: Learning where to sample model hypotheses. In *ICCV*, 2019. 2
- [5] Tat-Jun Chin, Pulak Purkait, Anders P. Eriksson, and David Suter. Efficient globally optimal consensus maximisation with tree search. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2017. 2
- [6] Sungjoon Choi, Qian-Yi Zhou, and Vladlen Koltun. Robust reconstruction of indoor scenes. In *CVPR*, 2015. 1, 2
- [7] Christopher Choy, JunYoung Gwak, and Silvio Savarese. 4D spatio-temporal ConvNets: Minkowski convolutional neural networks. In *CVPR*, 2019. 2, 3
- [8] Christopher Choy and Junha Lee. Open universal correspondence network. <https://github.com/chrischoy/open-ucn>, 2019. 8
- [9] Christopher Choy, Jaesik Park, and Vladlen Koltun. Fully convolutional geometric features. In *ICCV*, 2019. 2, 5, 6
- [10] Christopher B Choy, JunYoung Gwak, Silvio Savarese, and Manmohan Chandraker. Universal correspondence network. In *Advances in Neural Information Processing Systems*, 2016. 8
- [11] Ondrej Chum and Jiri Matas. Matching with prozac-progressive sample consensus. In *CVPR*, 2005. 1, 2
- [12] Angela Dai, Matthias Nießner, Michael Zollöfer, Shahram Izadi, and Christian Theobalt. BundleFusion: Real-time globally consistent 3d reconstruction using on-the-fly surface re-integration. *ACM Transactions on Graphics*, 2017. 2
- [13] Zheng Dang, Kwang Moo Yi, Yinlin Hu, Fei Wang, Pascal Fua, and Mathieu Salzmann. Eigendecomposition-free training of deep networks with zero eigenvalue-based losses. In *ECCV*, 2018. 2
- [14] Wei Dong, Jaesik Park, Yi Yang, and Michael Kaess. GPU accelerated robust scene reconstruction. In *IROS*, 2019. 2
- [15] Martin A. Fischler and Robert C. Bolles. Random sample consensus: A paradigm for model fitting with applications to image analysis and automated cartography. *Communications of the ACM*, 1981. 2, 6
- [16] Andrew W. Fitzgibbon. Robust registration of 2D and 3D point sets. *Image and Vision Computing*, 2003. 1, 2
- [17] Ben Glocker, Shahram Izadi, Jamie Shotton, and Antonio Criminisi. Real-time RGB-D camera relocalization. In *ISMAR*, 2013. 5
- [18] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 1, 2
- [19] Jared Heinly, Johannes L Schonberger, Enrique Dunn, and Jan-Michael Frahm. Reconstructing the world in six days. In *CVPR*, 2015. 7
- [20] M. B. Horowitz, N. Matni, and J. W. Burdick. Convex relaxations of SE(2) and SE(3) for visual pose estimation. In *ICRA*, 2014. 2
- [21] Reza Hoseinnezhad and Alireza Bab-Hadiashar. An estimator for high breakdown robust estimation in computer vision. *Computer Vision and Image Understanding*, 2011. 1, 2
- [22] Peter J Huber. *Robust statistics*. Springer, 2011. 1, 2
- [23] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*, 2015. 1
- [24] Gregory Izatt and Russ Tedrake. Globally optimal object pose estimation in point clouds with mixed-integer programming. In *International Symposium on Robotics Research*, 2017. 2
- [25] Karel Lebeda, Jiri Matas, and Ondrej Chum. Fixing the locally optimized RANSAC—full experimental evaluation. In *British Machine Vision Conference*, 2012. 2

Table 4: Classification scores on the YFCC100M test set. We consider a correspondence a true positive if the symmetric epipolar distance is below 10^{-4} and the prediction confidence is above 0.5.

	Yi <i>et al.</i> [48]			Zhang <i>et al.</i> [51]			Ours			Ours + SC			Ours + UCN		
	Prec.	Recall	F1	Prec.	Recall	F1	Prec.	Recall	F1	Prec.	Recall	F1	Prec.	Recall	F1
BUCKINGHAM	0.497	0.772	0.605	0.486	0.889	0.629	0.535	0.822	0.648	0.611	0.835	0.705	0.769	0.892	0.826
NOTRE DAME	0.581	0.894	0.705	0.629	0.951	0.757	0.647	0.915	0.758	0.721	0.929	0.812	0.844	0.954	0.896
REICHTAG	0.747	0.877	0.807	0.734	0.917	0.815	0.695	0.911	0.789	0.769	0.897	0.827	0.856	0.937	0.895
SACRE COEUR	0.658	0.871	0.750	0.662	0.948	0.780	0.632	0.917	0.748	0.718	0.932	0.811	0.855	0.948	0.899
Average	0.621	0.854	0.717	0.628	0.926	0.745	0.628	0.891	0.736	0.704	0.898	0.789	0.830	0.933	0.879

- [26] Hongdong Li. Consensus set maximization with guaranteed global optimality for robust geometry estimation. In *ICCV*, 2009. 2
- [27] David G. Lowe. Distinctive image features from scale-invariant keypoints. *International Journal of Computer Vision*, 2004. 7
- [28] Haggai Maron, Nadav Dym, Itay Kezurer, Shahar Kovalsky, and Yaron Lipman. Point registration via efficient convex relaxation. *ACM Transactions on Graphics*, 2016. 2
- [29] OpenMP Architecture Review Board. OpenMP application program interface version 3.0, May 2008. 3
- [30] G. Dias Pais, Pedro Miraldo, Srikumar Ramalingam, Venu Madhav Govindu, Jacinto C. Nascimento, and Rama Chellappa. 3DRegNet: A deep neural network for 3D point registration. *arXiv*, 2019. 2
- [31] Jaesik Park, Qian-Yi Zhou, and Vladlen Koltun. Colored point cloud registration revisited. In *ICCV*, 2017. 2
- [32] Thomas Probst, Danda Pani Paudel, Ajad Chhatkuli, and Luc Van Gool. Unsupervised learning of consensus maximization for 3D vision problems. In *CVPR*, 2019. 2
- [33] Charles R. Qi, Hao Su, Kaichun Mo, and Leonidas J. Guibas. Pointnet: Deep learning on point sets for 3D classification and segmentation. In *CVPR*, 2017. 1, 2, 3, 4
- [34] Rahul Raguram, Ondrej Chum, Marc Pollefeys, Jiri Matas, and Jan-Michael Frahm. USAC: A universal framework for random sample consensus. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2013. 1, 2
- [35] René Ranftl and Vladlen Koltun. Deep fundamental matrix estimation. In *ECCV*, 2018. 1, 2
- [36] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *MICCAI*, 2015. 1, 2
- [37] David M Rosen, Luca Carlone, Afonso S Bandeira, and John J Leonard. Se-sync: A certifiably correct algorithm for synchronization over the special euclidean group. *International Journal of Robotics Research*, 2019. 2
- [38] Peter J. Rousseeuw. Least median of squares regression. *Journal of the American Statistical Association*, 1984. 1, 2
- [39] Radu Bogdan Rusu, Nico Blodow, and Michael Beetz. Fast point feature histograms (FPFH) for 3d registration. In *ICRA*, 2009. 1, 5, 6, 7
- [40] Schönberger, Johannes Lutz, and Jan-Michael Frahm. Structure-from-motion revisited. In *CVPR*, 2016. 1
- [41] Ruwan B. Tennakoon, Alireza Bab-Hadiashar, Zhenwei Cao, Reza Hoseinnezhad, and David Suter. Robust model fitting using higher than minimal subset sampling. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2016. 1, 2
- [42] Bart Thomee, David A Shamma, Gerald Friedland, Benjamin Elizalde, Karl Ni, Douglas Poland, Damian Borth, and Li-Jia Li. YFCC100M: The new data in multimedia research. *Communications of the ACM*, 2016. 7
- [43] Philip HS Torr and Andrew Zisserman. MLESAC: A new robust estimator with application to estimating image geometry. *Computer Vision and Image Understanding*, 2000. 2
- [44] Philip H. S. Torr. Bayesian model estimation and selection for epipolar geometry and generic manifold fitting. *International Journal of Computer Vision*, 2002. 2
- [45] Jianxiong Xiao, Andrew Owens, and Antonio Torralba. Sun3D: A database of big spaces reconstructed using sfm and object labels. In *ICCV*, 2013. 5
- [46] Jiaolong Yang, Hongdong Li, and Yunde Jia. Go-ICP: Solving 3D registration efficiently and globally optimally. In *ICCV*, 2013. 2
- [47] Jiaolong Yang, Hongdong Li, and Yunde Jia. Optimal essential matrix estimation via inlier-set maximization. In *ECCV*, 2014. 2
- [48] Kwang Moo Yi, Eduard Trulls, Yuki Ono, Vincent Lepetit, Mathieu Salzmann, and Pascal Fua. Learning to find good correspondences. In *CVPR*, 2018. 1, 2, 3, 4, 5, 6, 8, 9, 10
- [49] Manzil Zaheer, Satwik Kottur, Siamak Ravanbakhsh, Barnabas Poczos, Ruslan R Salakhutdinov, and Alexander J Smola. Deep sets. In *Advances in Neural Information Processing Systems*. 2017. 1, 2, 3, 4
- [50] Andy Zeng, Shuran Song, Matthias Nießner, Matthew Fisher, Jianxiong Xiao, and Thomas Funkhouser. 3DMatch: Learning local geometric descriptors from RGB-D reconstructions. In *CVPR*, 2017. 5, 6
- [51] Jiahui Zhang, Dawei Sun, Zixin Luo, Anbang Yao, Lei Zhou, Tianwei Shen, Yurong Chen, Long Quan, and Hongen Liao. Learning two-view correspondences and geometry using order-aware network. In *ICCV*, 2019. 1, 2, 7, 8, 9, 10
- [52] Zhengyou Zhang. Determining the epipolar geometry and its uncertainty: A review. *International Journal of Computer Vision*, 1998. 2
- [53] Qian-Yi Zhou, Jaesik Park, and Vladlen Koltun. Fast global registration. In *ECCV*, 2016. 1, 2, 6
- [54] Qian-Yi Zhou, Jaesik Park, and Vladlen Koltun. Open3D: A modern library for 3D data processing. *arXiv*, 2018. 1, 6

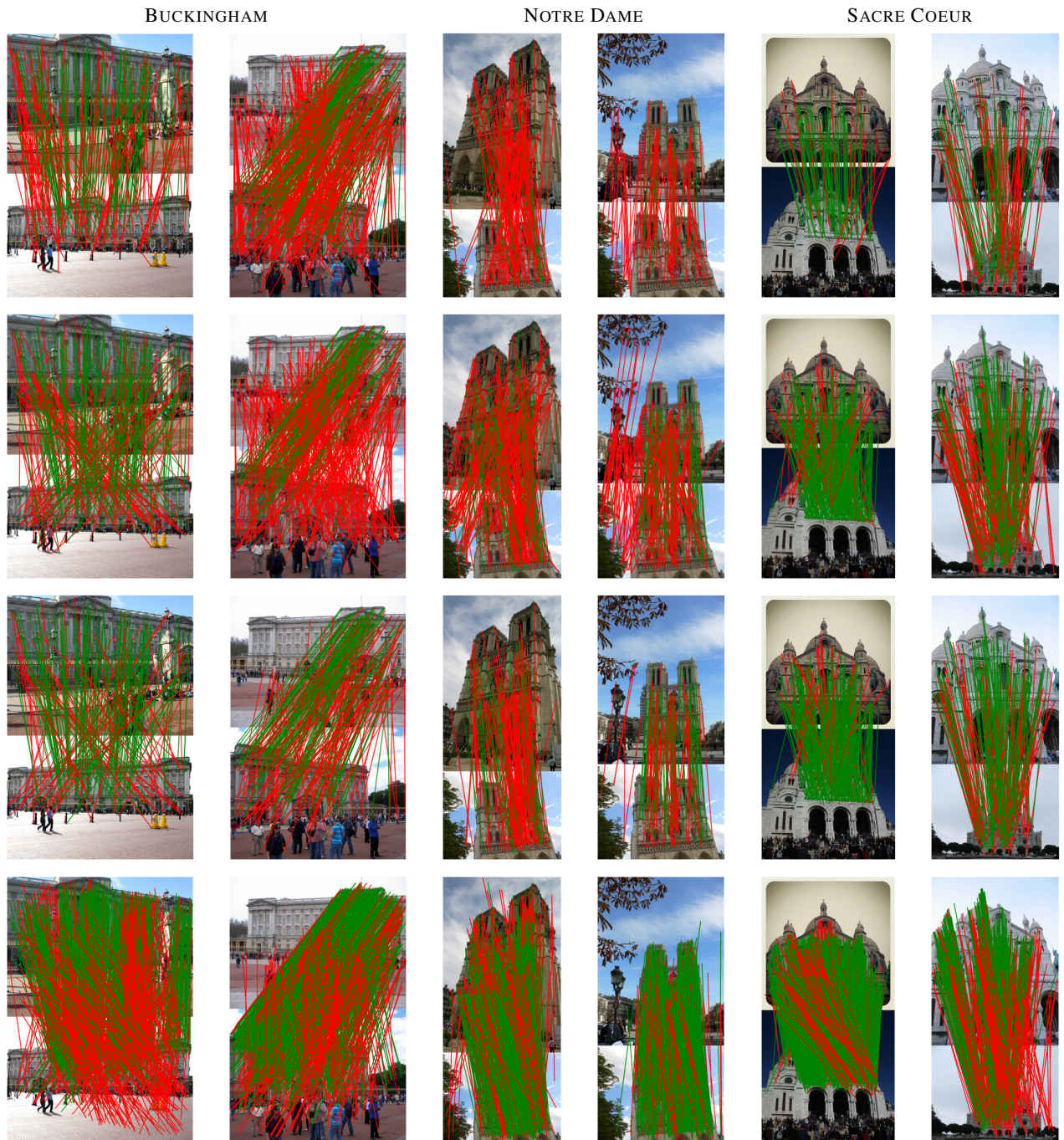


Figure 7: Matching results using Yi *et al.* [48] (first row), Zhang *et al.* [51] (second row), Ours + SC (third row) and Ours + UCN (last row). We visualize correspondences with inlier probability above 0.5. We color a correspondence green if it is a true positive (symmetric epipolar distance smaller than 10^{-4}) and red otherwise.