



HAL
open science

PCAM: Product of Cross-Attention Matrices for Rigid Registration of Point Clouds

Anh-Quan Cao, Gilles Puy, Alexandre Boulch, Renaud Marlet

► **To cite this version:**

Anh-Quan Cao, Gilles Puy, Alexandre Boulch, Renaud Marlet. PCAM: Product of Cross-Attention Matrices for Rigid Registration of Point Clouds. ICCV 2021 - International Conference on Computer Vision, Oct 2021, Montreal / Virtual, Canada. 10.1109/iccv48922.2021.01298 . hal-03364975

HAL Id: hal-03364975

<https://hal.science/hal-03364975v1>

Submitted on 5 Oct 2021

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

PCAM: Product of Cross-Attention Matrices for Rigid Registration of Point Clouds

Anh-Quan Cao^{2,*}

Gilles Puy¹

Alexandre Boulch¹

Renaud Marlet^{1,3}

¹Valeo.ai, Paris, France

²Inria, Paris, France[†]

³LIGM, Ecole des Ponts, Univ Gustave Eiffel, CNRS, Marne-la-Vallée, France

Abstract

Rigid registration of point clouds with partial overlaps is a longstanding problem usually solved in two steps: (a) finding correspondences between the point clouds; (b) filtering these correspondences to keep only the most reliable ones to estimate the transformation. Recently, several deep nets have been proposed to solve these steps jointly. We built upon these works and propose PCAM: a neural network whose key element is a pointwise product of cross-attention matrices that permits to mix both low-level geometric and high-level contextual information to find point correspondences. These cross-attention matrices also permits the exchange of context information between the point clouds, at each layer, allowing the network construct better matching features within the overlapping regions. The experiments show that PCAM achieves state-of-the-art results among methods which, like us, solve steps (a) and (b) jointly via deepnets. Our code and trained models are available at <https://github.com/valeoai/PCAM>.

1. Introduction

Point cloud registration is the problem of estimating the rigid transformation that aligns two point clouds. It has many applications in various domains such as autonomous driving, motion and pose estimation, 3D reconstruction, simultaneous localisation and mapping (SLAM), and augmented reality. The most famous method to solve this task is ICP [3], for which several improvements have been proposed, modifying the original optimisation process [42, 47] or using geometric feature descriptors [33] to match points.

Recently, end-to-end learning-based methods combining point feature extraction, point matching, and point-pairs filtering, have been developed to solve this task. Deep Closest Point (DCP) [40], improved by PRNet [39], finds point cor-

respondences via an attention matrix and estimate the transformation by solving a least-squares problem. Deep Global Registration (DGR) [7] tackles partial-to-partial point cloud registration by finding corresponding points via deep features, computing confidence scores for these pairs, and solving a weighted least-squares problem to estimate the transformation. IDAM [23] considers all possible pairs of points in a similarity matrix computed using deep or hand-crafted features, and proposes a learned two-step filter to select only relevant pairs for transformation estimation. Finally, [9, 12, 43] use deep networks and the Sinkhorn algorithm to find and filter corresponding points.

Our observation is that, when matching points between point clouds, one would like to find correspondences using both local fine geometric information, to precisely select the best corresponding point, and high-level contextual information, to differentiate between points with similar local geometry but from different parts of the scene. The fine geometric information is naturally extracted at the first layers of a deep convnet, while the context information is found in the deepest layers. The existing deep registration methods estimate the correspondences between point clouds using features extracted at the deepest layers. Therefore, there is no explicit control on the amount of fine geometric and high-level context information encoded in these features.

Instead, we propose to compute point correspondences at every layer of our deep network via cross-attention matrices, and to combine these matrices via a point-wise multiplication. This simple yet very effective solution naturally ensures that both low-level geometric and high-level context information are exploited when matching points. It also permits to remove spurious matches found only at one scale. Furthermore, we also exploit these cross-attention matrices to exchange information between the point clouds at each layer, allowing the network to exploit context information from both point clouds to find the best matching point within the overlapping regions.

The design of our method is inspired by DGR [7], DCP [40] and PRNet [39]. We exploit cross-attention matrices to

*Most of the work was done during an internship at valeo.ai in 2020.

[†]Inria, Mines ParisTech, PSL Research University.

exchange context information between point clouds, like in [40, 39, 12], to find the best matches within overlapping region. Our main contribution is to propose to compute such cross-attention matrices at every network layer and to combine them to exploit both fine and high-level information when matching points. Our second contribution is to show that our method achieves state-of-the-art results on two real datasets (indoor, outdoor) and on a synthetic one.

2. Related Work

Optimisation-based methods. Iterative Closest Point (ICP) [6, 3, 46] is the most known algorithm for point cloud registration. It takes in two point clouds and alternate between point matching via nearest-neighbour search and transformation estimation by solving a least-squares problem. Several improvements of ICP’s steps have been introduced to improve speed, solve the discretisation problem, become robust to outliers, or incorporate confidence scores on the correspondences [1, 4, 29, 31, 34]. ICP often converges to a local minimum due to the non-convexity of the objective function. Therefore, some works propose solutions to widen the basin of attraction of the global optimum [13, 37], use genetic algorithms [35], or find a good crude alignment to initialise ICP [26]. Another line of works concerns algorithms with theoretical convergence guarantees to a global optimum thanks to, *e.g.*, the combination of ICP and a branch-and-bound technique [42], via convex relaxations [27, 30], or thanks to mixed-integer programming [19]. While theoretically appealing, these approaches usually suffers from a high computational complexity.

Feature-based point matching. Instead of relying solely on point coordinates to find correspondences between points, several approaches have proposed to extract point feature descriptors by analysing the local, and possibly global, 3D geometry. Classical methods use hand-crafted features such as [20, 32, 33, 36] and, more recently, features extracted thanks to deep networks that take as inputs either hand-crafted local features [10, 11, 16, 18, 21, 45], or directly point cloud coordinates [8, 18]. These methods use these point features to establish correspondences and rely on RANSAC to estimate the transformation, possibly using a pre-filtering of the points unlikely to provide good correspondences [18].

End-to-end learning. Another type of approaches consists in training a network that will establish correspondences between the points of both point clouds, filter these correspondences to keep only the reliable pairs, and estimate the transformation based on the best pairs of points. Several methods fall in this category, such as [7, 9, 12, 23, 39, 40, 43] which we described in the introduction. In addition, let us also mention PointNetLK [2] that extracts a global feature vector per point cloud and unroll the Lucas-Kanade algorithm [25] to find the best transforma-

tion aligning the point clouds in feature space, DeepGMR [44] that maps each point cloud to a parametric probability distribution via a neural network and estimates the transformation by minimising the KL-divergence between distributions. [22] proposes a probabilistic registration method using deep features and learned attention weights. Finally, [15] proposes a method for multiview registration by jointly training a pairwise registration module and a global refinement module.

3. Network Architecture

3.1. Problem Statement

We consider the problem of partial-to-partial rigid registration between two point clouds $\mathcal{P} = \{\mathbf{p}_1, \dots, \mathbf{p}_N\}$ and $\mathcal{Q} = \{\mathbf{q}_1, \dots, \mathbf{q}_M\}$, with $\mathbf{p}_i, \mathbf{q}_j \in \mathbb{R}^3$. The point clouds \mathcal{P} and \mathcal{Q} represent two views of the same scene with partial overlap. Hence, only a subset of the points in \mathcal{P} can have matching points in \mathcal{Q} , and vice versa from \mathcal{Q} to \mathcal{P} . Let \mathcal{P}_v (resp. \mathcal{Q}_v) be the subset of points in \mathcal{P} (resp. \mathcal{Q}) visible in \mathcal{Q} (resp. \mathcal{P}). Our goal is to estimate the rotation matrix $\mathbf{R}_{\text{gt}} \in \text{SO}(3)$ and translation vector $\mathbf{t}_{\text{gt}} \in \mathbb{R}^3$ that aligns \mathcal{P}_v on \mathcal{Q}_v . This transformation can be estimated by solving

$$\min_{\mathbf{R}, \mathbf{t}} \sum_{\mathbf{p} \in \mathcal{P}_v} \|\mathbf{R}\mathbf{p} + \mathbf{t} - m_{\mathcal{Q}}(\mathbf{p})\|_2^2, \quad (1)$$

where $m_{\mathcal{Q}} : \mathcal{P}_v \rightarrow \mathcal{Q}_v$ maps any $\mathbf{p} \in \mathcal{P}_v$ to its best matching point $\mathbf{q} \in \mathcal{Q}_v$. Note that we can also estimate the inverse transformation from \mathcal{Q} to \mathcal{P} with the map $m_{\mathcal{P}} : \mathcal{Q}_v \rightarrow \mathcal{P}_v$.

Problem (1) can be solved via an SVD decomposition [3, 17]. The challenging tasks are the estimations of the subset of points \mathcal{P}_v and \mathcal{Q}_v , and of the mappings $m_{\mathcal{Q}}(\cdot)$ and $m_{\mathcal{P}}(\cdot)$, given only \mathcal{P} and \mathcal{Q} as inputs.

3.2. Method Overview

Our method is composed of two classical modules: point matching and point-pair filtering (see Figure 1). The point-matching module, denoted by $g(\cdot, \cdot)$, permits us to estimate two maps $\tilde{m}_{\mathcal{Q}} : \mathcal{P} \rightarrow \mathcal{Q}$ and $\tilde{m}_{\mathcal{P}} : \mathcal{Q} \rightarrow \mathcal{P}$ that provide pairs of corresponding points between \mathcal{P} and \mathcal{Q} (even in non-overlapping regions). The point-pair filtering module, denoted by $h(\cdot, \cdot)$, provides a confidence score to all pairs $(\mathbf{p}_i, \tilde{m}_{\mathcal{Q}}(\mathbf{p}_i))_{1 \leq i \leq N}$, $(\mathbf{q}_j, \tilde{m}_{\mathcal{P}}(\mathbf{q}_j))_{1 \leq j \leq M}$, hence allowing detection of valid pairs in overlapping regions.

Our main contribution is in g , where we propose the construction of cross-attention matrices at each layer of g and their combination to obtain point correspondences. These successive cross-attention matrices are obtained using point features with an increasing field of view, or scale, thus establishing point correspondences with increasing scene context. We obtain $\tilde{m}_{\mathcal{Q}}$ and $\tilde{m}_{\mathcal{P}}$ by combining all these attention matrices via pointwise matrix multiplications. Within overlapping regions, it permits us to improve the quality of

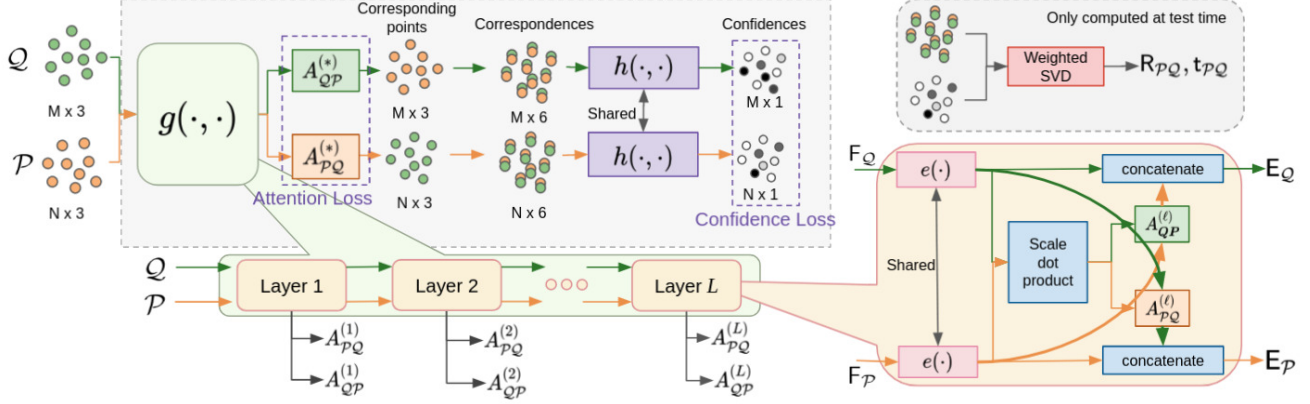


Figure 1. **Overview of our architecture.** \mathcal{P} and \mathcal{Q} enter in the point matching module $g(\cdot, \cdot)$ that permits the extraction of attention matrices used to compute one matching point for each input point. Each pair of matching points is then given a confidence score via $h(\cdot, \cdot)$, which permits the estimation of the rigid transformation by solving a weighted least-squares problem via SVD.

the maps \tilde{m}_Q and \tilde{m}_P by filtering out the matches that are inconsistent across scales. We also use these cross-attention matrices to exchange information between \mathcal{P} and \mathcal{Q} : each point in \mathcal{P} receives, in addition to its own feature, the feature of the best corresponding point in \mathcal{Q} (and vice-versa) before computing the point feature at the next layer. This process favours the discovery or rejection of correspondences at the next layer depending on the level of similarity between corresponding points at the current layer.

3.3. First Module: Point Matching

Our point matching module, denoted by $g: (\mathbb{R}^{N \times 3}, \mathbb{R}^{M \times 3}) \rightarrow (\mathbb{R}^{N \times M}, \mathbb{R}^{M \times N})$, takes two unaligned point clouds \mathcal{P} and \mathcal{Q} as input, and produces two global attention matrices $A_{PQ}^{(*)}$ and $A_{QP}^{(*)}$ which inform us on the correspondences between point clouds. These global attention matrices are constructed (see Sec. 3.3.2) from L attention matrices, $A_{PQ}^{(\ell)}$ and $A_{QP}^{(\ell)}$, obtained at layers $\ell = 1, \dots, L$ of our deep network (see Sec. 3.3.1).

3.3.1 Layer-wise Cross-Attention Matrices $A_{PQ}^{(\ell)}, A_{QP}^{(\ell)}$

Let $F_P \in \mathbb{R}^{N \times c^{(\ell-1)}}$ and $F_Q \in \mathbb{R}^{M \times c^{(\ell-1)}}$ be the features on \mathcal{P} and \mathcal{Q} , respectively, at the input of layer ℓ , where $c^{(\ell-1)}$ is the size of the feature vectors at the output of layer $\ell - 1$. Each point is thus described by a $c^{(\ell-1)}$ -dimensional feature vector. These features enter the ℓ^{th} point convnet $e: \mathbb{R}^{N \times c^{(\ell-1)}} \rightarrow \mathbb{R}^{N \times c^{(\ell)}/2}$ to extract new features of size $c^{(\ell)}/2$ for each input point, to be doubled after concatenation of the best corresponding point feature in the other point cloud (see below). This encoder extracts local geometry information around each point in the point cloud thanks to three residual blocks, each containing two FKACnv [5] sub-layers. The supplementary material details $e(\cdot)$.

The new feature vectors $e(F_P)$ and $e(F_Q)$ are used to compute two attention matrices $A_{PQ}^{(\ell)}, A_{QP}^{(\ell)} \in \mathbb{R}^{N \times M}$ with

$$(A_{PQ}^{(\ell)})_{ij} = \frac{e^{a_{ij}/s}}{\sum_{k=1}^M e^{a_{ik}/s}}, \quad (A_{QP}^{(\ell)})_{ij} = \frac{e^{a_{ij}/s}}{\sum_{k=1}^N e^{a_{kj}/s}}, \quad (2)$$

where $s > 0$ is the softmax temperature and

$$a_{ij} = \frac{e(F_P)_i e(F_Q)_j^T}{\|e(F_P)_i\|_2 \|e(F_Q)_j\|_2}. \quad (3)$$

A_{PQ}, A_{QP} differ in the normalisation dimension: A_{PQ} transfers information from \mathcal{Q} to \mathcal{P} , A_{QP} from \mathcal{P} to \mathcal{Q} .

Finally, we transfer information between point clouds by computing $E_P = [e(F_P), A_{PQ}^{(\ell)} e(F_Q)] \in \mathbb{R}^{N \times c^{(\ell)}}$, $E_Q = [e(F_Q), A_{QP}^{(\ell)T} e(F_P)] \in \mathbb{R}^{M \times c^{(\ell)}}$. E_P, E_Q are the output features of layer ℓ and input features of layer $(\ell + 1)$. At the end of the process, we have extracted two sets of L attention matrices: $(A_{PQ}^{(1)}, \dots, A_{PQ}^{(L)})$ and $(A_{QP}^{(1)}, \dots, A_{QP}^{(L)})$.

3.3.2 Global Attention Matrices $A_{PQ}^{(*)}, A_{QP}^{(*)}$

We combine the attention matrices of each layer ℓ via a simple pointwise multiplication, denoted by \odot , to obtain

$$A_{PQ}^{(*)} = A_{PQ}^{(1)} \odot \dots \odot A_{PQ}^{(L)}. \quad (4)$$

$A_{QP}^{(*)}$ is defined similarly. The motivation for this strategy is that the successive attention matrices are constructed using features with an increasing field of view or scale, and we want to match points only if their features are similar at *all* scales, hence the entry-wise multiplication of the attention coefficients. Another motivation is to permit to backpropagate gradients directly from the loss to each layer ℓ .

3.3.3 Soft and Sparse Maps

There exists two standards ways to match points once the global attention matrices are available: soft and sparse mappings. We explore the performance of both approaches in this work. The soft map $\tilde{m}_{\mathcal{Q}}(\cdot)$ is defined as

$$\tilde{m}_{\mathcal{Q}}(\mathbf{p}_i) = \frac{\sum_{j=1}^M (A_{\mathcal{P}\mathcal{Q}}^{(*)})_{i,j} \mathbf{q}_j}{\sum_{k=1}^M (A_{\mathcal{P}\mathcal{Q}}^{(*)})_{i,k}} \quad (5)$$

for points $(\mathbf{p}_i)_{1 \leq i \leq N}$. It maps \mathcal{P} to \mathbb{R}^3 rather than \mathcal{Q} . The sparse map $\tilde{m}_{\mathcal{Q}}(\cdot)$ is defined as

$$\tilde{m}_{\mathcal{Q}}(\mathbf{p}_i) = \mathbf{q}_{j^*}, \text{ where } j^* = \underset{j}{\operatorname{argmax}} (A_{\mathcal{P}\mathcal{Q}}^{(*)})_{ij}. \quad (6)$$

It maps \mathcal{P} to \mathcal{Q} , but it is not differentiable and it prevents backpropagation from the confidence network to the point-matching network. Yet, the point-matching network will still be trained using a cross-entropy loss on the attention matrices (Sec. 3.6). The mapping $\tilde{m}_{\mathcal{P}}$ from \mathcal{Q} to \mathcal{P} or \mathbb{R}^3 is constructed similarly using $A_{\mathcal{Q}\mathcal{P}}^{(*)}$.

3.4. Second Module: Confidence Estimation

The module in Sec.3.3 produces pairs of matching points $(\mathbf{p}_i, \tilde{m}_{\mathcal{Q}}(\mathbf{p}_i))$ and $(\mathbf{q}_j, \tilde{m}_{\mathcal{P}}(\mathbf{q}_j))$ for all points in \mathcal{P} and \mathcal{Q} . However, as we are tackling partial-to-partial registration, only subsets of \mathcal{P} and \mathcal{Q} match to one another. Hence, there are many incorrect matches which need to be filtered out. We detect these incorrect matches by concatenating each pair in a vector $[\mathbf{p}_i^T, \tilde{m}_{\mathcal{Q}}(\mathbf{p}_i)^T]^T \in \mathbb{R}^6$, and then pass these N pairs into the confidence estimation module $h : \mathbb{R}^{\times 6} \rightarrow (0, 1)^{\times 1}$ which outputs a confidence score $w_{\mathbf{p}_i}$ for each input pair. The same module h is used for pairs $(\mathbf{q}_j, \tilde{m}_{\mathcal{P}}(\mathbf{q}_j))$ and yields M corresponding scores $w_{\mathbf{q}_j}$. The confidence estimator is a point convnet with 9 residual blocks, each containing 2 FKACnv sub-layers. The confidence score are thus estimated using context information provided by the point convolution on \mathcal{P} (or \mathcal{Q}). The detailed network architecture is in the supplementary material.

3.5. Transformation Estimation

We estimate the rotation and translation from \mathcal{P} to \mathcal{Q} by solving a weighted least-squares problem, using the weights $w_{\mathbf{p}_i}$ given by the confidence estimator:

$$(\mathbf{R}_{\text{est}}, \mathbf{t}_{\text{est}}) = \underset{\mathbf{R}, \mathbf{t}}{\operatorname{argmin}} \sum_{i=1}^N \phi(w_{\mathbf{p}_i}) \|\mathbf{R}\mathbf{p}_i + \mathbf{t} - \tilde{m}_{\mathcal{Q}}(\mathbf{p}_i)\|_2^2. \quad (7)$$

The function $\phi : \mathbb{R} \rightarrow \mathbb{R}$ applies hard-thresholding on the confidence scores by setting to zero the less confident ones. The above problem can be solved with an SVD [17, 40].

3.6. Losses

The networks $g(\cdot, \cdot)$ and $h(\cdot, \cdot)$ and trained under supervision using a first loss, denoted by L^{ca} that applies on the attention matrices $A_{\mathcal{P}\mathcal{Q}}^{(*)}, A_{\mathcal{Q}\mathcal{P}}^{(*)}$ and the sum of other losses, L^{cc} and L^{gc} , that apply on the confidence scores $w_{\mathbf{p}_i}, w_{\mathbf{q}_j}$. The complete training loss satisfies $L^{\text{ca}} + L^{\text{cc}} + L^{\text{gc}}$.

3.6.1 Loss on the Attention Matrices

The role of the attention matrices $A_{\mathcal{P}\mathcal{Q}}^{(*)}, A_{\mathcal{Q}\mathcal{P}}^{(*)}$ is to identify mappings between \mathcal{P} and \mathcal{Q} for the final registration task. The pairs of points identified at this stage should include, as a subset, the ideal pairs encoded by $m_{\mathcal{P}}$ and $m_{\mathcal{Q}}$. Identifying this subset is the role of the confidence estimator $h(\cdot, \cdot)$.

The ideal map $m_{\mathcal{P}}$ and $m_{\mathcal{Q}}$ are estimated as follows. For each point \mathbf{p}_i , we search the closest point $\mathbf{q}_{j^*(i)}$ to $\mathbf{R}_{\text{gt}}\mathbf{p}_i + \mathbf{t}_{\text{gt}}$ in \mathcal{Q} and consider the pair $(\mathbf{p}_i, \mathbf{q}_{j^*(i)})$ as valid if $\mathbf{R}_{\text{gt}}\mathbf{p}_i + \mathbf{t}_{\text{gt}}$ is also the closest point to $\mathbf{q}_{j^*(i)}$ in $\{\mathbf{R}_{\text{gt}}\mathbf{p}_u + \mathbf{t}_{\text{gt}}\}_u$. More precisely, let us define $j^*(u) = \underset{j}{\operatorname{argmin}}_j \|\mathbf{R}_{\text{gt}}\mathbf{p}_u + \mathbf{t}_{\text{gt}} - \mathbf{q}_j\|_2$ and $i^*(v) = \underset{u}{\operatorname{argmin}}_u \|\mathbf{R}_{\text{gt}}\mathbf{p}_i + \mathbf{t}_{\text{gt}} - \mathbf{q}_v\|_2$. The ideal maps satisfy $m_{\mathcal{Q}}(\mathbf{p}_u) = \mathbf{q}_v$, $m_{\mathcal{P}}(\mathbf{q}_v) = \mathbf{p}_u$, $\forall (u, v) \in \mathcal{C}$, where $\mathcal{C} = \{(u, v) \mid j^*(u) = v, u = i^*(v)\}$. Only a subset of the points are in \mathcal{C} , even in overlapping regions. For convenience, we also define the set of points in \mathcal{P} for which a corresponding point exists in \mathcal{Q} by $\mathcal{C}_{\mathcal{P}} = \{u \mid \exists v \in \llbracket M \rrbracket \text{ s.t. } (u, v) \in \mathcal{C}\}$. The set $\mathcal{C}_{\mathcal{Q}}$ is defined similarly for the inverse mapping.

We consider a contrastive classification loss, such as used in [7, 23, 9], denoted by $L^{\text{ca}} = L_{\mathcal{P}\mathcal{Q}}^{\text{ca}} + L_{\mathcal{Q}\mathcal{P}}^{\text{ca}}$. $L_{\mathcal{P}\mathcal{Q}}^{\text{ca}}$ enforces a good mapping from \mathcal{P} to \mathcal{Q} and satisfies

$$\begin{aligned} L_{\mathcal{P}\mathcal{Q}}^{\text{ca}} &= -\frac{1}{N} \sum_{(u,v) \in \mathcal{C}} \log \left[(A_{\mathcal{P}\mathcal{Q}}^{(*)})_{uv} \right] \\ &= -\frac{1}{N} \sum_{\ell=1}^L \sum_{(u,v) \in \mathcal{C}} \log \left[\frac{e^{(a_{uv}^{(\ell)})/s}}{\sum_{j=1}^M e^{(a_{uj}^{(\ell)})/s}} \right]. \end{aligned} \quad (8)$$

$L_{\mathcal{Q}\mathcal{P}}^{\text{ca}}$ is defined likewise. All points are constrained in (8) thanks to the softmax involved in the attention matrices. In cases where \mathcal{C} contains only few points, rather than augmenting \mathcal{C} with pairs of worse quality, we leave \mathcal{C} untouched to force the point matching network to identify the most reliable pairs (along with bad ones) and let the second module learn how to filter the bad pairs.

3.6.2 Losses on the Confidence Scores

To train the confidence estimator $h(\cdot, \cdot)$, we consider classification and geometric losses.

The classification losses are built by measuring the quality of the maps $\tilde{m}_{\mathcal{Q}}, \tilde{m}_{\mathcal{P}}$ at each input point. If a mapped point $\tilde{m}_{\mathcal{Q}}(\mathbf{p}_i)$ is close to the ideal target point $\mathbf{R}_{\text{gt}}\mathbf{p}_i + \mathbf{t}_{\text{gt}}$, then the pair $(\mathbf{p}_i, \tilde{m}_{\mathcal{Q}}(\mathbf{p}_i))$ is considered accurate and $w_{\mathbf{p}_i}$

should be close to 1, and 0 otherwise. The classification loss thus satisfies $L^{cc} = L_{\mathcal{P}\mathcal{Q}}^{cc} + L_{\mathcal{Q}\mathcal{P}}^{cc}$, where

$$L_{\mathcal{P}\mathcal{Q}}^{cc} = -\left(\frac{1}{N} \sum_i y_{\mathbf{p}_i} \log(w_{\mathbf{p}_i}) + (1 - y_{\mathbf{p}_i}) \log(1 - w_{\mathbf{p}_i})\right),$$

$$y_{\mathbf{p}_i} = \begin{cases} 1, & \text{if } \|\mathbf{R}_{\text{gt}}\mathbf{p}_i + \mathbf{t}_{\text{gt}} - \tilde{m}_{\mathcal{Q}}(\mathbf{p}_i)\|_2 \leq \kappa, \\ 0, & \text{otherwise,} \end{cases} \quad (9)$$

and $\kappa > 0$ is a threshold to decide if a pair of points is accurate or not. The second loss $L_{\mathcal{Q}\mathcal{P}}^{cc}$ is defined similarly by considering the inverse mapping.

The threshold $\kappa > 0$ to decide pair accuracy is somewhat arbitrary. To mitigate the effect of this choice, we also consider the geometric loss $L^{gc} = L_{\mathcal{P}\mathcal{Q}}^{gc} + L_{\mathcal{Q}\mathcal{P}}^{gc}$ where

$$L_{\mathcal{P}\mathcal{Q}}^{gc} = \sum_i \frac{w_{\mathbf{p}_i}}{N} \|\mathbf{R}_{\text{gt}}\mathbf{p}_i + \mathbf{t}_{\text{gt}} - \tilde{m}_{\mathcal{Q}}(\mathbf{p}_i)\|_2. \quad (10)$$

$L_{\mathcal{Q}\mathcal{P}}^{gc}$ is defined using $w_{\mathbf{q}_j}$ and the inverse transformation. Whenever the distance between a mapped point and its ideal location is large, then $w_{\mathbf{p}_i}$ should be small. On the contrary, when this distance is small, then $w_{\mathbf{p}_i}$ can be large. Note that this geometric losses can be used only in combination with the classification losses L^{cc} as otherwise $w_{\mathbf{p}_i}, w_{\mathbf{q}_j} = 0$ is a trivial useless solution.

4. Experiments

4.1. Datasets and Training Parameters

We evaluate our method on real (indoor, outdoor) and synthetic datasets. The indoor dataset is 3DMatch [45]. We use the standard train/test splits and the procedure of [7, 10, 11, 8] to generate pairs of scans with at least 30% of overlap for training and testing. During training, as in [7], we apply data augmentation using random rotations in $[0^\circ, 360^\circ)$ around a random axis, and random scalings in $[0.8, 1.2]$. For the experiments on outdoor data, we use the KITTI odometry dataset [14] and follow the same protocol as [7]: GPS-IMU is used to create pairs of scans that are at least 10m apart; the ground-truth transformation is computed using GPS followed by ICP. Unlike in [7], we do not use data augmentation during training on KITTI. For synthetic data, we use ModelNet40 [41] and follow the setup of [39] to simulate partial registration problems.

All models are trained using AdamW [24], with a weight decay of 0.001, a batch size of 1, and a learning rate of 0.001. On 3Dmatch and KITTI, the models are trained for 100 epochs with a learning rate divided by 10 after 60 and 80 epochs. On ModelNet40, it is sufficient to train the models for 10 epochs (with a learning rate divided by 10 after 6 and 8 epochs) to observe convergence. All results are reported using the models obtained at the last epoch. The temperature s in (2) is set to $s = 0.03$.

	ϕ	Opt.	Saf.	TE _{all}	RE _{all}	Recall	TE	RE
FGR [47]						42.7	0.11	4.08
RANSAC [32]						74.9	0.09	2.92
DCP [40]						3.2	0.21	8.42
PointNetLK [2]						1.6	0.21	8.04
DGR [7]				0.47	17.4	73.9	0.09	2.97
DGR [7]	✓			0.38	13.4	81.5	0.08	2.56
DGR [7]	✓	✓		0.33	11.8	86.6	0.07	2.34
DGR [7]	✓	✓	✓	0.25	9.5	91.2	0.07	2.42
PCAM-Sparse				0.44	16.3	74.9	0.08	2.98
PCAM-Sparse	✓			0.35	13.0	83.8	0.07	2.43
PCAM-Sparse	✓	✓		0.32	11.8	87.0	0.07	2.11
PCAM-Sparse	✓	✓	✓	0.23	8.9	92.4	0.07	2.16
PCAM-Soft				0.46	17.5	74.7	0.09	3.01
PCAM-Soft	✓			0.37	13.1	81.4	0.08	2.50
PCAM-Soft	✓	✓		0.33	11.9	85.6	0.07	2.12
PCAM-Soft	✓	✓	✓	0.24	9.8	91.3	0.07	2.25

Table 1. **Results on 3DMatch test split.** The recall is computed using 0.3 m and 15° as thresholds. The scores of all variants of DGR are obtained using the official implementation of DGR [7], and those of the other concurrent methods are reported from [7]. The columns ‘ ϕ ’, ‘Opt.’, and ‘Saf.’ indicate respectively whether hard-thresholding on the confidence scores, the pose refinement of [7], and the safeguard registration of [7], are used or not.

4.2. Metrics

The performance on 3DMatch and KITTI is measured using the metrics of [7]: the translation error (TE) and rotation error (RE) are defined as $\text{TE}(\mathbf{t}_{\text{est}}) = \|\mathbf{t}_{\text{est}} - \mathbf{t}_{\text{gt}}\|_2$ and $\text{RE}(\mathbf{R}_{\text{est}}) = \arccos[(\text{Tr}(\mathbf{R}_{\text{gt}}^T \mathbf{R}_{\text{est}}) - 1)/2]$, respectively. We also computed the metric coined ‘recall’ in [7], which is the percentage of registrations whose rotation and translation errors are both smaller than predefined thresholds, i.e., the proportion of successful registrations. We report the average rotation error RE_{all} and translation error TE_{all} on *all* pairs of scans, as well as average rotation error RE and translation error TE on the subset of *successful* registrations. Performance on ModelNet40 is computed using the metrics of [39].

4.3. Comparison with Existing Methods

We report in this section the performance of PCAM on 3DMatch and KITTI. The results obtained on ModelNet40 are available in the supplementary material and show that PCAM outperforms the concurrent methods on this dataset. In what follows, PCAM-Soft and PCAM-Sparse refer to our method with soft (5) and sparse maps (6), respectively.

	ϕ	Opt.	ICP	TE _{all}	RE _{all}	Recall	TE	RE
FGR [47]						0.2	0.41	1.02
RANSAC [32]						34.2	0.26	1.39
FCGF [8]						98.2	0.1	0.33
DGR [7]				0.77	2.32	63.1	0.28	0.56
DGR [7]	✓			0.76	2.32	63.4	0.28	0.56
DGR [7]	✓	✓		0.34	1.62	96.6	0.21	0.33
DGR [7]	✓	✓	✓	0.16	1.43	98.2	0.03	0.14
PCAM - Soft				0.18	1.00	97.2	0.08	0.33
PCAM - Sparse				0.22	1.17	96.5	0.08	0.31
PCAM - Soft			✓	0.12	0.79	98.0	0.03	0.14
PCAM - Sparse			✓	0.17	1.04	97.4	0.03	0.14

Table 2. **Results on KITTI test split.** The recall is computed using 0.6 m and 5° as thresholds. The scores of all variants of DGR are obtained using the official implementation of DGR [7], and those of the other concurrent methods are reported from [7]. The columns ‘ ϕ ’, ‘Opt.’, and ‘ICP’ indicate respectively whether hard-thresholding on the confidence scores, the pose refinement of [7], and post-processing by ICP, are used or not.

4.3.1 Indoor Dataset: 3DMatch

We train PCAM on 3DMatch using $L = 6$ layers, with point clouds obtained by sampling $N = M = 4096$ points at random from the voxelised point clouds (voxel size of 5 cm). The parameter κ in (9) is set at 12 cm.

For a thorough comparison with DGR [7], we study the performance of PCAM after applying each DGR’s post-processing steps: filtering the confidence weights with ϕ , refining the pose estimation using the optimisation proposed in [7], and using the same safeguard as in [7]. The threshold applied in ϕ is tuned on 3DMatch’s validation set. The safeguard is activated when the ℓ_1 -norm of the confidence weights is below a threshold, in which case the transformation is estimated using RANSAC (see details in [7]). For a fair comparison with DGR, we compute this second threshold such that DGR and PCAM use the safeguard on the same proportion of scans.

We report the performance of PCAM and concurrent methods in Table 1. PCAM achieves the best results and this is confirmed in the curves of Fig. 2. We provide illustrations of the quality of matched points and registrations in the supplementary material. We did not compare PCAM with RLL, which is not performing as well as DGR on 3DMatch [22].

4.3.2 Outdoor Dataset: KITTI

We train PCAM on 3DMatch using $L = 6$ layers with point clouds of (at most) $N = M = 2048$ points drawn at random from the voxelised point clouds (voxel size of 30 cm). The parameter κ in (9) is set at 60 cm.

The performance of our method is compared to others

in Table 2. As scores after refinement of DGR registration by ICP are reported on this dataset in [7], we also include the scores of PCAM after the same refinement in this table. PCAM outperforms DGR on all metrics with both sparse (6) and soft (5) maps, except on the recall when using sparse maps, where it is just 0.1 point behind DGR. When combined with ICP, PCAM also achieves better results than DGR except on the recall but where it is just 0.2 point behind with soft maps and 0.8 with sparse maps. Finally, we notice in Fig. 2 a clear advantage of PCAM over DGR for the estimation of the translation. Note that PCAM achieves better performance than DGR without the need to use ϕ , nor used the pose refinement of [7]; we did not notice any significant improvement when using them on KITTI. We compare PCAM to RLL on the version of KITTI used in [22]: PCAM outperforms RLL with a recall of 84.7% for sparse maps and 86.5% for soft maps, vs 76.9% for RLL.

4.4. Ablation Study

In this section, we conduct two ablation studies: the first justifies our choice of the training loss; the second demonstrates the benefit of the product of attention matrices and of the exchange of contextual information between point clouds. The performance of the trained models are evaluated on the validation set of the considered datasets. On 3DMatch and KITTI, the input point clouds are obtained by drawing at random (at most) $N = M = 2048$ points from the voxelised point clouds. Because of these random draws, different evaluations lead to different scores. To take into account these variations, we report the average scores and the standard deviations obtained using three evaluations on the validation set for each trained model. On ModelNet40, we do not subsample the point clouds. The rotation R_{est} and translation t_{est} are computed using the method described in Sec. 3.5 *without* using any hard-thresholding on the confidence scores. We conduct experiments with $L = 2$ or $L = 6$ layers. For the experiments at $L = 2$, two models are trained from different initialisation and the reported average scores and the standard deviations are computed using both models. One model is trained for the experiments at $L = 6$.

4.4.1 Training Loss

The experiments in this section are conducted with $L = 2$ layers and using the product of attention matrices (4).

We recall that the loss L^{cc} cannot be removed as there exists a trivial solution where all the confidence scores are equal to zero in absence of this loss. Hence, the only loss that can be removed is L^{gc} . The results in Table 3 show that removing L^{gc} yields worse performance both with soft and sparse mappings. We explain this result by the fact that R_{est} and t_{est} are obtained using a slightly modified version of L^{gc} at test time (see Sec. 3.5), hence the presence of L^{gc} at

					3DMatch			KITTI			ModelNet40			
L^{ca}	L^{ga}	L^{gc}	$\odot_{\ell} A^{(\ell)}$	$A^{(L)}$	A	L	Rec. (%)	RE _{all}	TE _{all}	Rec. (%)	RE _{all}	TE _{all}	RMSE (R) ($\times 10^{-3}$)	RMSE (t) ($\times 10^{-3}$)
Soft map	✓	✓	✓	✓		2	53.7 (1.8)	36.4 (1.9)	0.49 (0.03)	93.4 (0.6)	3.0 (0.2)	0.45 (0.02)	18 (1.2)	0.13 (0.002)
	✓		✓	✓		2	52.3 (0.6)	36.3 (1.0)	0.49 (0.01)	94.3 (0.5)	2.5 (0.1)	0.42 (0.02)	15 (2.7)	0.09 (0.014)
	✓			✓		2	45.9 (1.5)	41.9 (2.5)	0.55 (0.03)	93.7 (0.3)	3.0 (0.1)	0.44 (0.01)	18 (1.1)	0.12 (0.015)
		✓	✓	✓		2	19.4 (2.7)	67.4 (1.5)	0.85 (0.02)	90.9 (1.1)	3.1 (0.1)	0.52 (0.04)	494 (218)	1.78 (0.879)
	✓		✓		✓	2	37.1 (2.9)	47.9 (2.2)	0.65 (0.02)	90.4 (0.9)	3.5 (0.1)	0.61 (0.08)	38 (3.9)	0.28 (0.007)
	✓		✓			✓ 2	37.3 (1.4)	47.2 (0.6)	0.65 (0.01)	92.7 (0.4)	2.8 (0.2)	0.53 (0.01)	66 (10.5)	0.39 (0.010)
	✓		✓	✓		6	82.4 (0.3)	13.5 (0.3)	0.21 (0.01)	95.4 (0.3)	2.2 (0.1)	0.33 (0.01)	16	0.11
	✓		✓		✓	6	76.6 (1.9)	18.2 (0.8)	0.27 (0.01)	93.5 (0.8)	2.9 (0.03)	0.46 (0.03)	22	0.16
	✓		✓			✓ 6	64.4 (1.1)	27.8 (0.6)	0.38 (0.01)	93.9 (0.0)	2.0 (0.1)	0.33 (0.01)	28	0.16
Sparse map	✓	N/A	✓	✓		2	52.4 (2.1)	38.6 (2.0)	0.51 (0.02)	94.6 (0.5)	2.8 (0.5)	0.39 (0.02)	19 (1.9)	0.10 (0.016)
	✓	N/A		✓		2	48.5 (1.0)	39.4 (1.5)	0.53 (0.02)	94.4 (0.4)	2.6 (0.4)	0.36 (0.03)	24 (6.6)	0.14 (0.031)
	✓	N/A	✓		✓	2	54.1 (3.8)	36.0 (2.7)	0.49 (0.03)	94.2 (0.5)	2.6 (0.2)	0.43 (0.05)	15 (0.4)	0.09 (0.002)
	✓	N/A	✓			✓ 2	50.9 (1.3)	38.9 (1.5)	0.52 (0.02)	94.6 (0.8)	2.7 (0.3)	0.43 (0.03)	17 (1.3)	0.09 (0.004)
	✓	N/A	✓	✓		6	81.6 (1.2)	16.3 (1.2)	0.24 (0.01)	95.9 (0.3)	2.3 (0.05)	0.33 (0.01)	36	0.20
	✓	N/A	✓		✓	6	78.4 (1.1)	17.5 (1.3)	0.24 (0.01)	96.5 (0.3)	2.4 (0.1)	0.46 (0.003)	27	0.15
	✓	N/A	✓			✓ 6	68.0 (1.1)	25.9 (0.5)	0.34 (0.01)	95.7 (0.3)	1.9 (0.1)	0.33 (0.01)	18	0.10

Table 3. **Ablation study on the validation set of 3DMatch, KITTI, and ModelNet40.** The losses L^{ca} , L^{gc} , L^{ga} are defined in the main paper. The global attention matrix $A^{(*)}$ is either equal to the pointwise multiplication of all the attention matrices ($\odot_{\ell} A^{(\ell)}$), or to the last attention matrix ($A^{(L)}$), or to the attention matrix (A) computed at the last layer of our network in a version where we have removed all the intermediate attention matrices $\ell = 1, \dots, L$ (see main paper for details).

training time probably improves the quality of confidence scores for the estimation of the transformation.

With soft maps (5), we can use another loss to train g in replacement or in complement to L^{ca} . The motivation for this loss is that L^{ca} penalises all points that are not mapped correctly to the ideal target point in the same way, whether the mapped point is close to the ideal location or far away. Instead, we can consider the possibility to use a geometric loss, $L^{ga} = L_{\mathcal{P}\mathcal{Q}}^{ga} + L_{\mathcal{Q}\mathcal{P}}^{ga}$, that takes into account the distance between the estimated and ideal corresponding points:

$$L_{\mathcal{P}\mathcal{Q}}^{ga} = \frac{1}{|\mathcal{C}_{\mathcal{P}}|} \sum_{u \in \mathcal{C}_{\mathcal{P}}} \|\tilde{m}_{\mathcal{Q}}(\mathbf{p}_u) - m_{\mathcal{Q}}(\mathbf{p}_u)\|_2. \quad (11)$$

$L_{\mathcal{Q}\mathcal{P}}^{ga}$ is defined likewise using $\mathcal{C}_{\mathcal{Q}}$. Note that it can be used to train g only in the case of soft maps (5) as $\tilde{m}_{\mathcal{Q}}$ is not differentiable when using sparse maps (6). The results in Table 3 show that the performance drops significantly when using L^{ga} alone. Note that a similar observation was made in [48]. Finally, using $L^{ga} + L^{ca}$ yields mildly better performance only on 3DMatch compared to using L^{ca} alone.

In view of these results, we chose to train PCAM using $L^{ca} + L^{cc} + L^{gc}$, both for soft and sparse maps.

4.4.2 Role of the Attention Matrices

Product of attention matrices: $\odot_{\ell} A^{(\ell)}$ vs $A^{(L)}$. To show the benefit of mixing low-level geometric information and high-level context information via the proposed product of attention matrices (4), we compare the performance of PCAM when replacing (4) by $A^{(*)} = A^{(L)}$, *i.e.*, using only the cross-attention matrix at the deepest layer. Note that the matrices $A^{(1)}, \dots, A^{(L-1)}$ are still present in g .

We expect the interest of the mix of information (4) to appear at large L (so that there is enough information to mix) and on challenging datasets where small overlaps and large transformations yield many ambiguities to be resolved for accurate point matching. In that respect, the most challenging dataset used in this work is 3DMatch. One can indeed verify in Table 3 that all models yields quite similar and very good results on KITTI and ModelNet40, leaving little room to clearly observe the impact of different network architectures. On ModelNet40, we highlight that all scores are close to zero (the scores are multiplied by 10^3 in Table 3) which makes it difficult to rank the different strategies as the difference between them could be due to training noise. In contrast, the impact of different choices of architecture is much more visible on 3DMatch.

We clearly observe the benefit of increasing L from 2

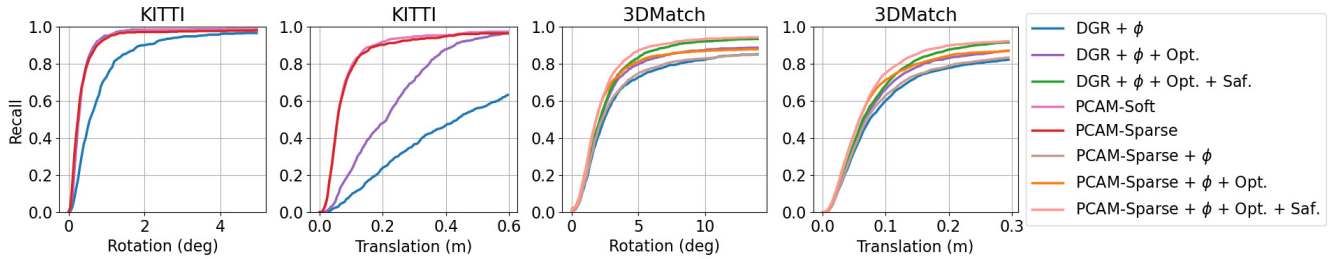


Figure 2. Registration recall vs rotation and translation error thresholds for KITTI (left) and 3DMatch (right).

to 6, and of mixing low-level geometric information and high-level context via the product of cross-attention matrices at $L = 6$ on 3DMatch. At $L = 2$, we observe a benefit of the product of attention matrices when using soft maps on 3DMatch, while this advantage is not visible when using sparse maps. We explain this fact because the product can make the soft maps more tightly concentrated around the found corresponding points, yielding better quality maps. This effect is invisible with sparse maps at $L = 2$ because, by construction, these maps are tightly concentrated around the corresponding points that are found.

Sharing contextual information: $A^{(L)}$ vs A . To show the benefit of exchanging information between point clouds via the intermediate cross-attention matrices $A^{(1)}, \dots, A^{(L-1)}$, we conduct experiments where we remove all intermediate matrices but the last one, which we denote by A . As before, we expect the benefit of this strategy to be more visible at large L and on challenging datasets such as 3DMatch. The results in Table 3 confirm this result where the benefit of exchanging contextual information is particularly noticeable at $L = 6$ on 3DMatch.

5. Discussion

In this section we differentiate our method from the most closely related works on three main aspects: exchange of contextual information between point clouds, computation of the pairs of corresponding points, and point-pair filtering.

Exchange of contextual information. Multiple cross-attention layers in PCAM allows the network to progressively exchange information between point clouds at *every* layer to find the best matches in overlapping regions. DCP [40], PRNet [39], and OPRNet [9] use transformer layers that also exchange information between point clouds, but only in the deepest layers of their convnet. PREDATOR [18] uses one cross-attention layer in the middle of their U-Net, hence only exchanging high-level contextual information. The newly published method [12] uses several cross-attention layers but, unlike us, does not merge them for point matching. There is no exchange of information between point clouds in [7, 23, 43].

Point matching. PCAM uses *multiple* cross-attention matrices that allows the network to exploit both fine-grained

geometric information and high-level context to find pairs of corresponding points. In contrast, DCP [40] and PRNet [39] use only one attention layer at the deepest layer of their network, which might limit their performance due to the lack of this fine-grained geometric information. RPM-Net[43], OPRNet [9] use a single Sinkhorn layer to establish the correspondences. IDAM [23] replaces the dot product operator by a learned convolution module that computes the similarity score between point features. The point-matching module in [15] is equivalent to ours in the non-optimal setup of one attention A .

Point filtering. PCAM first constructs pairs of corresponding points and assign to each of them a confidence score using a point convolution network. A similar strategy is used in [15]. DCP [40] does not have any filtering step, hence cannot handle partial-to-partial registration. Before starting to match points, PRNet [39] first selects a subset of points to be matched in each point cloud using the ℓ_2 -norm of point features. Similarly to PRNet [39], PREDATOR [18] also start with a selection of points to be matched in each point cloud, using the learned overlap probability and matchability scores. DGR [7] uses a similar process to PCAM but uses 6D convolutions to predict the confidence score, whereas we use 3D convolutions with 6D point-features (the concatenation of the coordinates of the point itself and its match point). IDAM [23] proposes a two-stage point elimination technique where points are individually filtered first; pairs of points are then discarded. Finally, [9, 12, 43] use their Sinkhorn layer to filter outliers.

6. Conclusion

In this work, we have proposed a novel network architecture where we transfer information between point clouds to find the best match within the overlapping regions, and mix low-level geometric information and high-level contextual knowledge to find correspondences between point clouds. Our experiments show that this architecture achieves state-of-the-art results on several datasets.

Acknowledgements. This work was partly performed using HPC resources from GENCI-IDRIS (Grant 2020-AD011012040). We thank Raoul de Charette for his constructive feedbacks on an earlier version of this paper.

References

- [1] G. Agamennoni, S. Fontana, R. Y. Siegwart, and D. G. Sorrenti. Point Clouds Registration with Probabilistic Data Association. In *International Conference on Intelligent Robots and Systems (IROS)*, pages 4092–4098, 2016.
- [2] Y. Aoki, H. Goforth, R. A. Srivatsan, and S. Lucey. PointNetLK: Robust & Efficient Point Cloud Registration Using PointNet. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7156–7165, 2019.
- [3] P. J. Besl and N. D. McKay. A method for registration of 3-D shapes. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 14(2):239–256, 1992.
- [4] S. Bouaziz, A. Tagliasacchi, and M. Pauly. Sparse Iterative Closest Point. *Computer Graphics Forum*, 32(5):113–123, 2013.
- [5] A. Boulch, G. Puy, and R. Marlet. FKACnv: Feature-Kernel Alignment for Point Cloud Convolution. In *Asian Conference on Computer Vision (ACCV)*, 2020.
- [6] Y. Chen and G. Medioni. Object modeling by registration of multiple range images. In *International Conference on Robotics and Automation (ICRA)*, volume 3, pages 2724–2729, 1991.
- [7] C. Choy, W. Dong, and V. Koltun. Deep Global Registration. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2511–2520, 2020.
- [8] C. Choy, J. Park, and V. Koltun. Fully convolutional geometric features. In *International Conference on Computer Vision (ICCV)*, pages 8957–8965, 2019.
- [9] Z. Dang, F. Wang, and M. Salzmann. Learning 3D-3D Correspondences for One-Shot Partial-to-partial Registration. *arXiv:2006.04523*, 2020.
- [10] H. Deng, T. Birdal, and S. Ilic. PPF-FoldNet: Unsupervised Learning of Rotation Invariant 3D Local Descriptors. In *European Conference on Computer Vision (ECCV)*, 2018.
- [11] H. Deng, T. Birdal, and S. Ilic. PPFNet: Global Context Aware Local Features for Robust 3D Point Matching. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 195–205, 2018.
- [12] K. Fischer, M. Simon, F. Olsner, S. Milz, H.-M. Gross, and P. Mader. StickyPillars: Robust and Efficient Feature Matching on Point Clouds Using Graph Neural Networks. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 313–323, 2021.
- [13] A. Fitzgibbon. Robust Registration of 2D and 3D Point Sets. *Image and Vision Computing*, 21:1145–1153, 2002.
- [14] 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.
- [15] Z. Gojcic, C. Zhou, J. Wegner, L. Guibas, and T. Birdal. Learning Multiview 3D Point Cloud Registration. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.
- [16] Z. Gojcic, C. Zhou, J. Wegner, and A. Wieser. The Perfect Match: 3D Point Cloud Matching With Smoothed Densities. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5540–5549, 2019.
- [17] J. Gower. Generalized procrustes analysis. *Psychometrika*, 40:33–51, 1975.
- [18] S. Huang, Z. Gojcic, M. Usvyatsov, A. Wieser, and K. Schindler. PREDATOR: Registration of 3D Point Clouds with Low Overlap. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4267–4276, 2021.
- [19] G. Izatt, H. Dai, and R. Tedrake. Globally Optimal Object Pose Estimation in Point Clouds with Mixed-Integer Programming. In *International Symposium on Robotics Research*, 2017.
- [20] A. E. Johnson and M. Hebert. Using spin images for efficient object recognition in cluttered 3d scenes. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 21(5):433–449, 1999.
- [21] M. Houry, Q.-Y. Zhou, and V. Koltun. Learning Compact Geometric Features. In *International Conference on Computer Vision (ICCV)*, pages 153–161, 2017.
- [22] F. J. Lawin and P.-E. Forssén. Registration loss learning for deep probabilistic point set registration. In *International Conference on 3D Vision (3DV)*, 2020.
- [23] J. Li, C. Zhang, Z. Xu, H. Zhou, and C. Zhang. Iterative Distance-Aware Similarity Matrix Convolution with Mutual-Supervised Point Elimination for Efficient Point Cloud Registration. In *European Conference on Computer Vision (ECCV)*, 2020.
- [24] I. Loshchilov and F. Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations (ICLR)*, 2019.
- [25] B. Lucas and T. Kanade. An Iterative Image Registration Technique with an Application to Stereo Vision. In *International joint conference on Artificial intelligence*, volume 2, pages 674–679, 1981.
- [26] A. Makadia, A. Patterson, and K. Daniilidis. Fully automatic registration of 3d point clouds. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, volume 1, pages 1297–1304, 2006.
- [27] H. Maron, N. Dym, I. Kezurer, S. Kovalsky, and Y. Lipman. Point Registration via Efficient Convex Relaxation. *ACM Transactions on Graphics*, 35(4):1–12, 2016.
- [28] V. Nair and G. E. Hinton. Rectified linear units improve restricted boltzmann machines. In *International Conference on Machine Learning (ICML)*, pages 807–814, 2010.
- [29] F. Pomerleau, F. Colas, and R. Siegwart. *A Review of Point Cloud Registration Algorithms for Mobile Robotics*. 2015.
- [30] D.M. Rosen, L. Carlone, A.S. Bandeira, and J.J. Leonard. A Certifiably Correct Algorithm for Synchronization over the Special Euclidean Group. In *International Workshop on the Algorithmic Foundations of Robotics (WAFR)*, volume 13, 2020.
- [31] S. Rusinkiewicz and M. Levoy. Efficient variants of the ICP algorithm. In *International Conference on 3-D Digital Imaging and Modeling*, pages 145–152, 2001.
- [32] R. B. Rusu, N. Blodow, and M. Beetz. Fast Point Feature Histograms (FPFH) for 3D registration. In *International Conference on Robotics and Automation (ICRA)*, pages 3212–3217, 2009.

- [33] S. Salti, F. Tombari, and L. Di Stefano. SHOT: Unique Signatures of Histograms for Surface and Texture Description. *Computer Vision and Image Understanding*, 125, 2014.
- [34] A. Segal, D. Hähnel, and S. Thrun. Generalized-ICP. In *Robotics: Science and Systems*, 2009.
- [35] L. Silva, O.R.P. Bellon, and K. L. Boyer. Precision range image registration using a robust surface interpenetration measure and enhanced genetic algorithms. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 27(5), 2005.
- [36] F. Tombari, S. Salti, and L. Di Stefano. Unique shape context for 3d data description. In *ACM workshop on 3D object retrieval*, 2010.
- [37] Y. Tsin and T. Kanade. A correlation-based approach to robust point set registration. In *European Conference on Computer Vision (ECCV)*, pages 558—569, 2004.
- [38] D. Ulyanov, A. Vedaldi, and V. Lempitsky. Instance normalization: The missing ingredient for fast stylization. *arXiv:1607.08022*, 2016.
- [39] Y. Wang and J. Solomon. PRNet: Self-Supervised Learning for Partial-to-Partial Registration. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 32, pages 8814–8826, 2019.
- [40] Y. Wang and J. M. Solomon. Deep Closest Point: Learning Representations for Point Cloud Registration. In *International Conference on Computer Vision (ICCV)*, pages 3522–3531, 2019.
- [41] Z. Wu, S. Song, A. Khosla, F. Yu, L. Zhang, X. Tang, and J. Xiao. 3D ShapeNets: A deep representation for volumetric shapes. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1912–1920, 2015.
- [42] J. Yang, H. Li, D. Campbell, and Y. Jia. Go-ICP: A Globally Optimal Solution to 3D ICP Point-Set Registration. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 38(11):2241–2254, 2016.
- [43] Z. J. Yew and G. H. Lee. Rpm-net: Robust point matching using learned features. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.
- [44] W. Yuan, B. Eckart, K. Kim, V. Jampani, D. Fox, and J. Kautz. Deepgmr: Learning latent gaussian mixture models for registration. In *European Conference on Computer Vision (ECCV)*, pages 733–750, 2020.
- [45] A. Zeng, S. Song, M. Nießner, M. Fisher, J. Xiao, and T. Funkhouser. 3DMatch: Learning Local Geometric Descriptors from RGB-D Reconstructions. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [46] Z. Zhang. Iterative point matching for registration of freeform curves and surfaces. *International Journal on Computer Vision*, 13(2):119–152, 1994.
- [47] Q.-Y. Zhou, J. Park, and V. Koltun. Fast Global Registration. In *European Conference on Computer Vision (ECCV)*, pages 766–782, 2016.
- [48] T. Zodage, R. Chakwate, V. Sarode, R. A. Srivatsan, and H. Choset. Correspondence matrices are underrated. In *International Conference on 3D Vision (3DV)*, 2020.

A. Network Architecture

A.1. Multiscale Attention Module $g(\cdot)$

Fig. 3 shows the detailed architecture of the multiscale attention module and of its L different encoders $e^{(\ell)} : \mathbb{R}^{\cdot \times c^{(\ell-1)}} \rightarrow \mathbb{R}^{\cdot \times c^{(\ell)}/2}$. We recall that $c^{(\ell)}$ denotes the number of channels at the output of the ℓ^{th} multiscale attention module. The encoder $e^{(\ell)}$ is made of three residual blocks, each yielding features with $c^{(\ell)}/2$ channels. The architecture of the residual block is also presented in Fig. 3 where the first 1D convolution is used only when the number of input and output channels differs. Each residual block consists of two FKACnv layers [5] with $c^{(\ell)}/2$ channels at input and output, a neighborhood of 32 nearest neighbours, a stride of size 1 (no downsampling of the point clouds), and a kernel of size 16. All convolutional layers are followed by instance normalization (IN) [38] with learned affine correction and a ReLU activation [28]. The number of channels $c^{(\ell)}$ used in our experiments is given in Table 4.

ℓ	0	1	2	3	4	5	6
$c^{(\ell)}$	3	32	32	64	64	128	128

Table 4. **Number of channels $c^{(\ell)}$ at the output of the ℓ^{th} block in the multiscale attention module.** $c^{(0)} = 3$ corresponds to the three coordinates x, y, z of the input point cloud.

A.2. Confidence Estimator $h(\cdot, \cdot)$

Fig. 3 shows the detailed architecture of the confidence estimator $h : \mathbb{R}^{\cdot \times 6} \rightarrow (0, 1)^{\cdot \times 1}$. It consists of nine residual blocks, a FKACnv layer that reduces the number of channels to 1 and a final sigmoid activation. The FKACnv layer has a neighborhood of 32 nearest neighbours, a stride of size 1, and a kernel of size 16. The residual blocks have the same structure as the one used to construct the encoders in the multiscale attention module.

B. Experimental results

B.1. Indoor dataset: 3DMatch

A detailed analysis of the performance of our method for each of the 8 scenes in the test set of 3DMatch [45] is presented in Fig. 4. The results confirm the global scores presented in the main part of the paper. Each of the post-processing steps permits to improve the recall on all the scenes. DGR and PCAM-Sparse achieve nearly similar recall on all scenes with a very noticeable difference on Hotel3 where PCAM-Sparse with hard-thresholding on the weights outperforms DGR with all post-processing steps used.

We present in Fig. 5 and Fig. 6 examples of success and failure, respectively, of PCAM-Sparse with filtering with ϕ

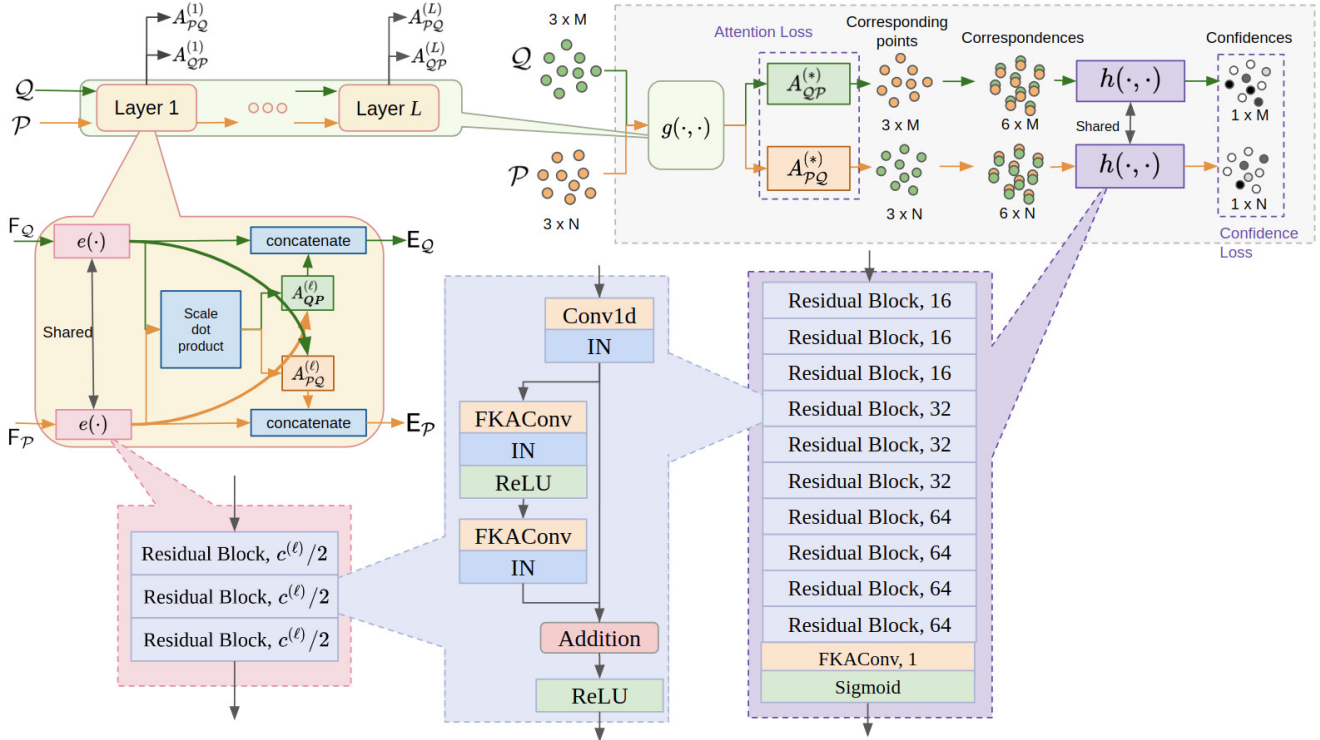


Figure 3. **Detailed network architecture of the multiscale attention module $g(\cdot, \cdot)$ and the confidence estimator module $h(\cdot, \cdot)$.** Each module is constructed using the same type residual block, in which all internal layers have the same number of channels C indicated in the notation “Residual block, C ”. Note that the 1D convolution in the residual block is used *only* when the number of input and output channels differ.

on 3DMatch [45]. We notice that the misregistrations are due to wrong mappings between points on planar surfaces in the examples of row 2, 3 in Fig. 6. We also remark in the example at row 5 that our method found matching points between points of similar but different object. Finally, row 1 of Fig. 6 shows a case where the estimated correspondences are correct but the confidence estimator considers them as unreliable which leads to a wrong estimation of the transformation (unless the safeguard is activated). This shows that the performance of the confidence estimator can still be improved.

B.2. Outdoor dataset: KITTI

We present in Fig. 7 and Fig. 8 examples of successful and failed registration results, respectively, with our method on the KITTI odometry dataset [14]. We observe in Fig. 7 that the static objects are well aligned after registrations. Concerning the failed registration results, it seems that this is due to mapping between similar but different structure in the scene, at least for the scans with very large registration errors such as on row 1 or 4.

B.3. Synthetic dataset: ModelNet40

We provide in Table 5 the results obtained for ModelNet40 on the split corresponding to unseen objects, unseen categories, and unseen objects with noise [39]. Our method achieves better results than the recent methods that provided scores on these versions of the dataset.

RPM-Net [43] is also evaluated on ModelNet40, but using a slightly different setting than the one used by the methods in Table 5. We also test PCAM on this variant of ModelNet40. RPM-Net uses several passes/iterations to align two point clouds while PCAM, which is not trained on small displacements for refinement, uses only one. For fairness, we evaluate both methods after the first main pass. PCAM outperforms RPM-Net on the ‘clean’ version of ModelNet40 (Chamfer error of $1.8 \cdot 10^{-5}$ for RPM-Net, $3.4 \cdot 10^{-9}$ for PCAM) and on its ‘noisy’ version ($7.9 \cdot 10^{-4}$ for RPM-Net, $6.9 \cdot 10^{-4}$ for PCAM).

We present in Fig. 9 examples of registration results with our method for pairs of scans in the unseen objects with Gaussian noise split. The results presented in these figures illustrate the accurate registrations as well as the good quality of the matched pairs of points.

Method	Unseen objects				Unseen categories				Unseen objects + noise			
	RMSE (R)	MAE (R)	RMSE (t)	MAE (t)	RMSE (R)	MAE (R)	RMSE (t)	MAE (t)	RMSE (R)	MAE (R)	RMSE (t)	MAE (t)
PointNetLK [2]	16.73	7.55	0.045	0.025	22.94	9.65	0.061	0.033	19.94	9.08	0.057	0.032
DCP-v2 [40]	6.71	4.45	0.027	0.020	9.77	6.95	0.034	0.025	6.88	4.53	0.028	0.021
PRNet [39]	3.20	1.45	0.016	0.010	4.98	2.33	0.021	0.015	4.32	2.05	0.017	0.012
IDAM [23]	2.46	0.56	0.016	0.003	3.04	0.61	0.019	0.004	3.72	1.85	0.023	0.011
OPRNet [9]	0.328	0.052	0.002	0.0003	0.357	0.069	0.002	0.0004	2.06	0.68	0.008	0.003
PCAM-Sparse	0.023	0.012	0.0002	0.00009	0.072	0.018	0.0002	0.0001	0.58	0.15	0.002	0.001
PCAM-Soft	0.023	0.006	0.0001	0.00004	0.056	0.013	0.0006	0.0001	0.47	0.15	0.002	0.001

Table 5. **Results on ModelNet40 for unseen objects, unseen categories and unseen objects with Gaussian noise.** The scores are reported from [39] for the first 3 methods, and from the associated papers for the others.

B.4. Sparse vs FKACnv convolutions

DGR uses sparse convolutions, whereas PCAM uses FKACnv point convolutions. To test if DGR suffers from a large disadvantage due to these sparse convolutions, we experiment replacing our point matching network by DGR’s pre-trained point matching network. We then retrain our confidence estimator on KITTI using DGR’s matched points. The performance of this new system reaches a recall of 94.6%, an RE_{all} of 2.9 and TE_{all} of 0.46 on KITTI validation set. This is on par with the results obtained with our point matching network, showing using sparse convolutions in the point matching network do not perform much worse than FKACnv.

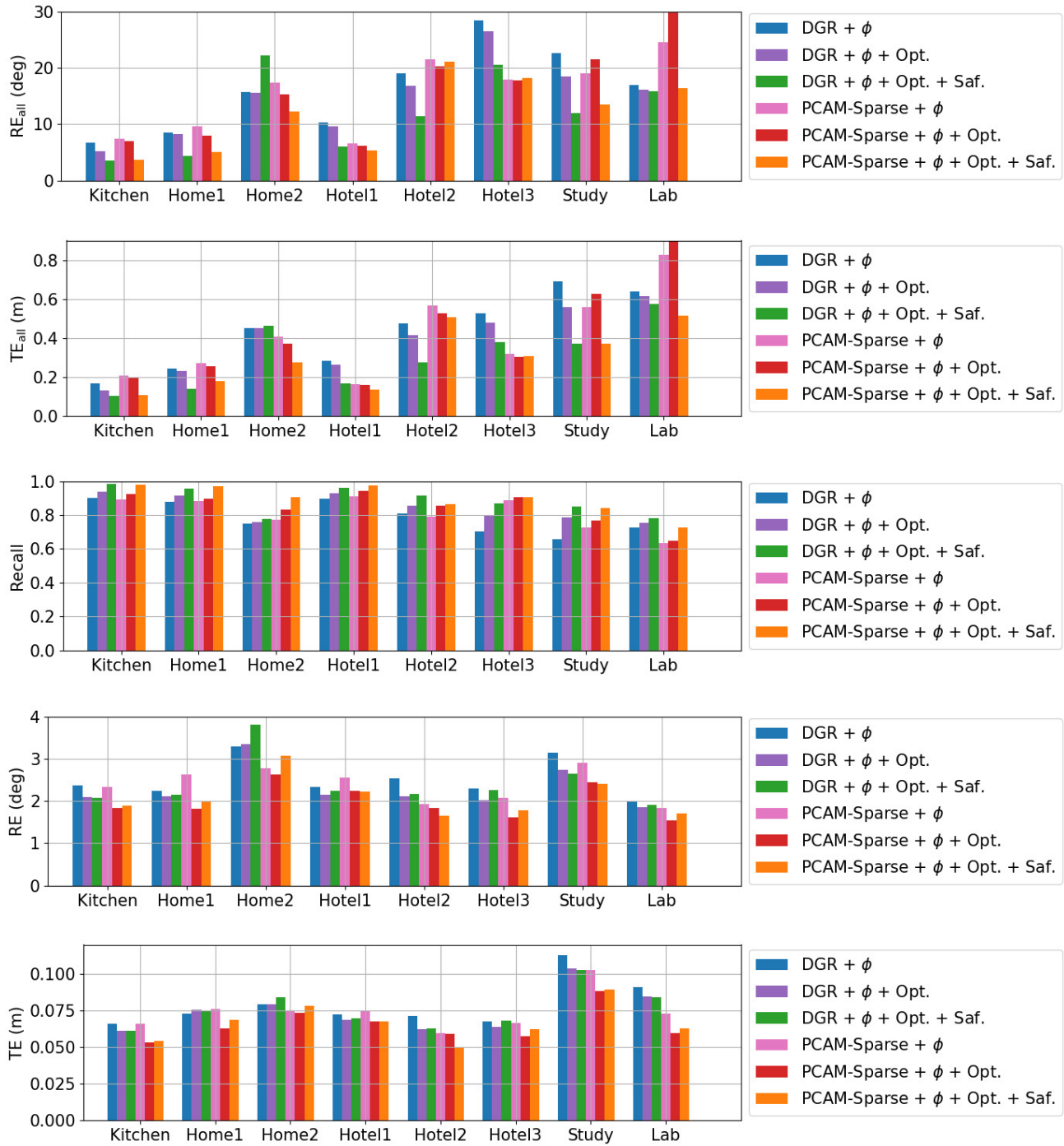


Figure 4. **Analysis of 3DMatch registration results per scene.** Row 1-2: Average TE and RE measured *on all pairs* (lower is better). Row 3: recall rate (higher is better). Row 4-5: TE and RE measured *on successfully registered pairs* (lower is better). Note that the recall of two methods have to be almost identical for their errors RE (resp. TE) to be comparable.

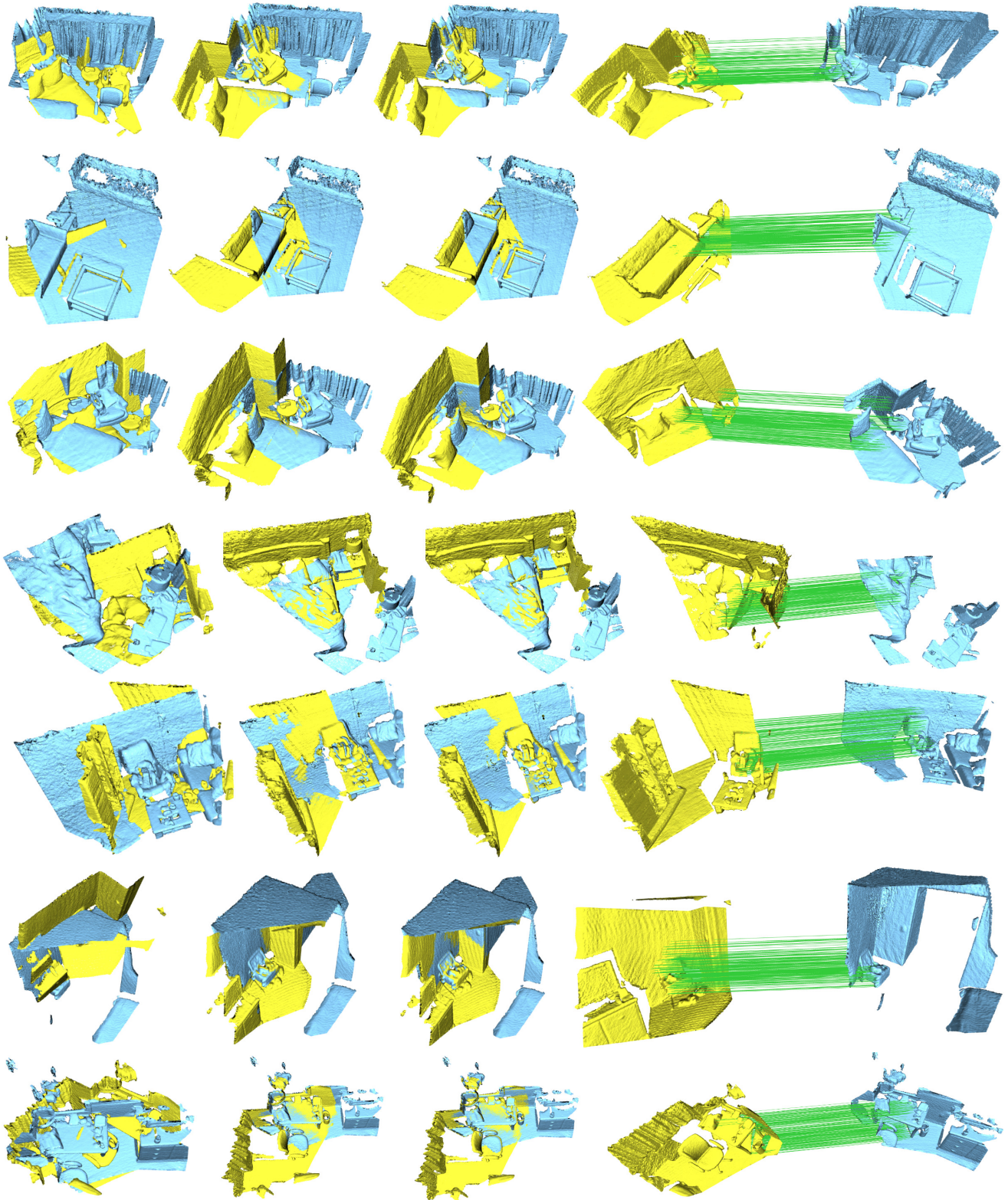


Figure 5. **Examples of successful registrations on 3DMatch with PCAM-Sparse + ϕ .** From left to right: overlaid, non-registered input scans (blue and yellow colors); ground-truth registration; scans registered with our method; and top 256 pairs of matched points with highest confidence.

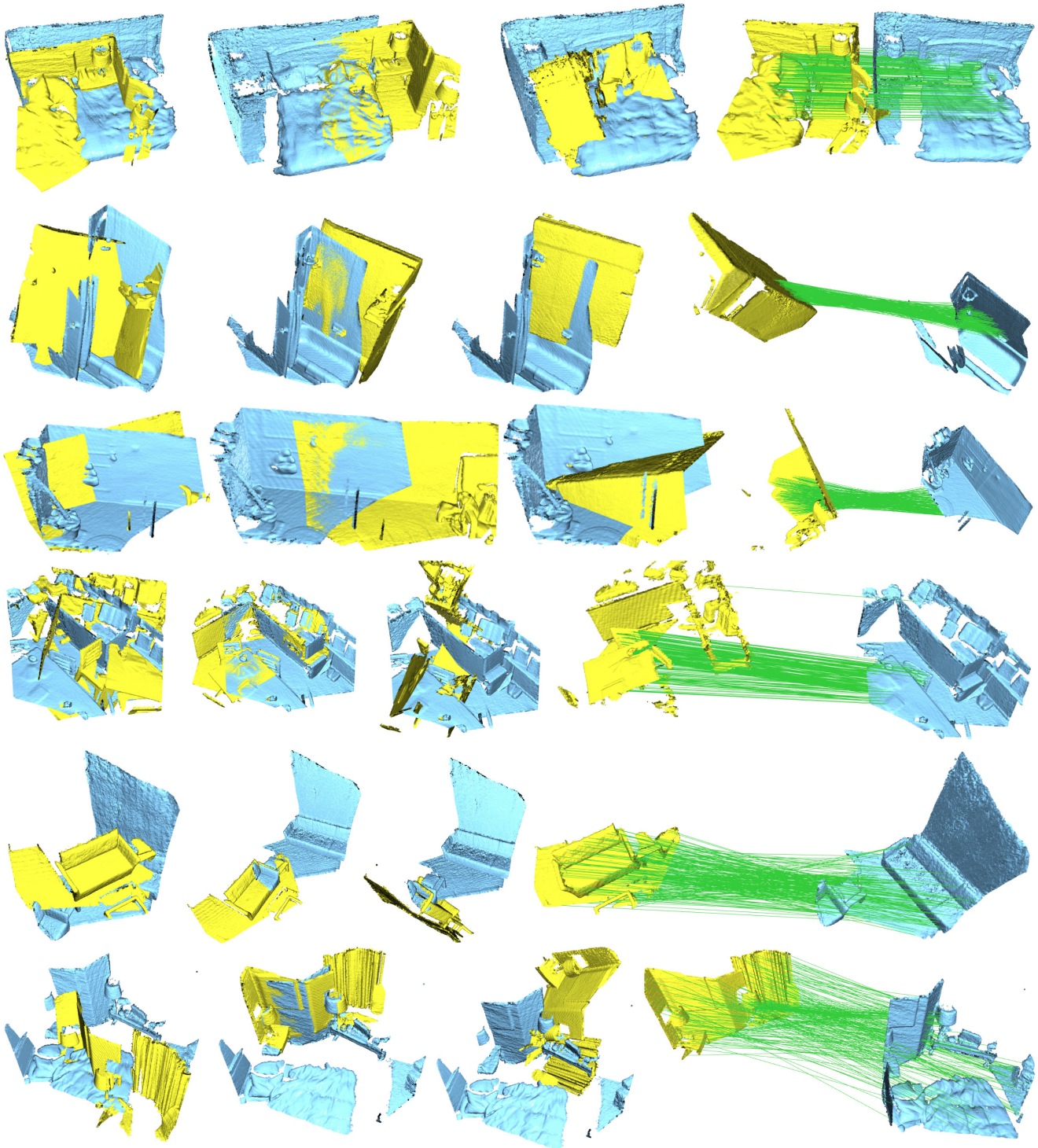


Figure 6. Examples of failed registrations on 3DMatch with PCAM-Sparse + ϕ . From left to right: overlaid, non-registered input scans (blue and yellow colors); ground-truth registration; scans registered with our method; and top 256 pairs of matched points with highest confidence.

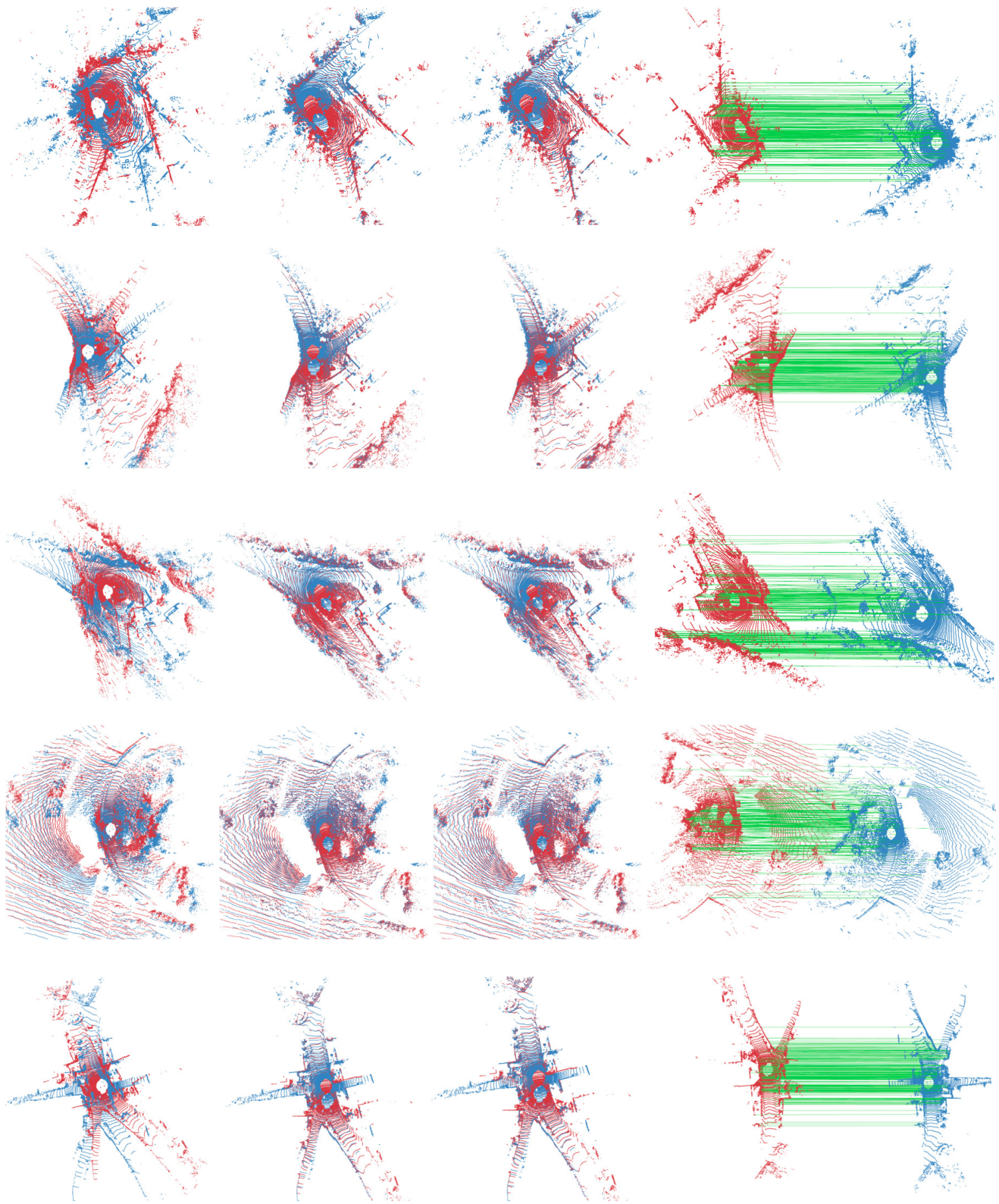


Figure 7. **Examples in bird-eye view of successful registrations on the KITTI odometry dataset [14].** From left to right: overlaid, non-registered input scans (blue and red colors); ground-truth registration; scans registered with our method; and top 256 pairs of matched points with highest confidence. *Note that we have used the full non-voxelized Lidar scans for a better visualisation. However, all registrations are done using 2048 points drawn at random after voxelization of the full point cloud.*

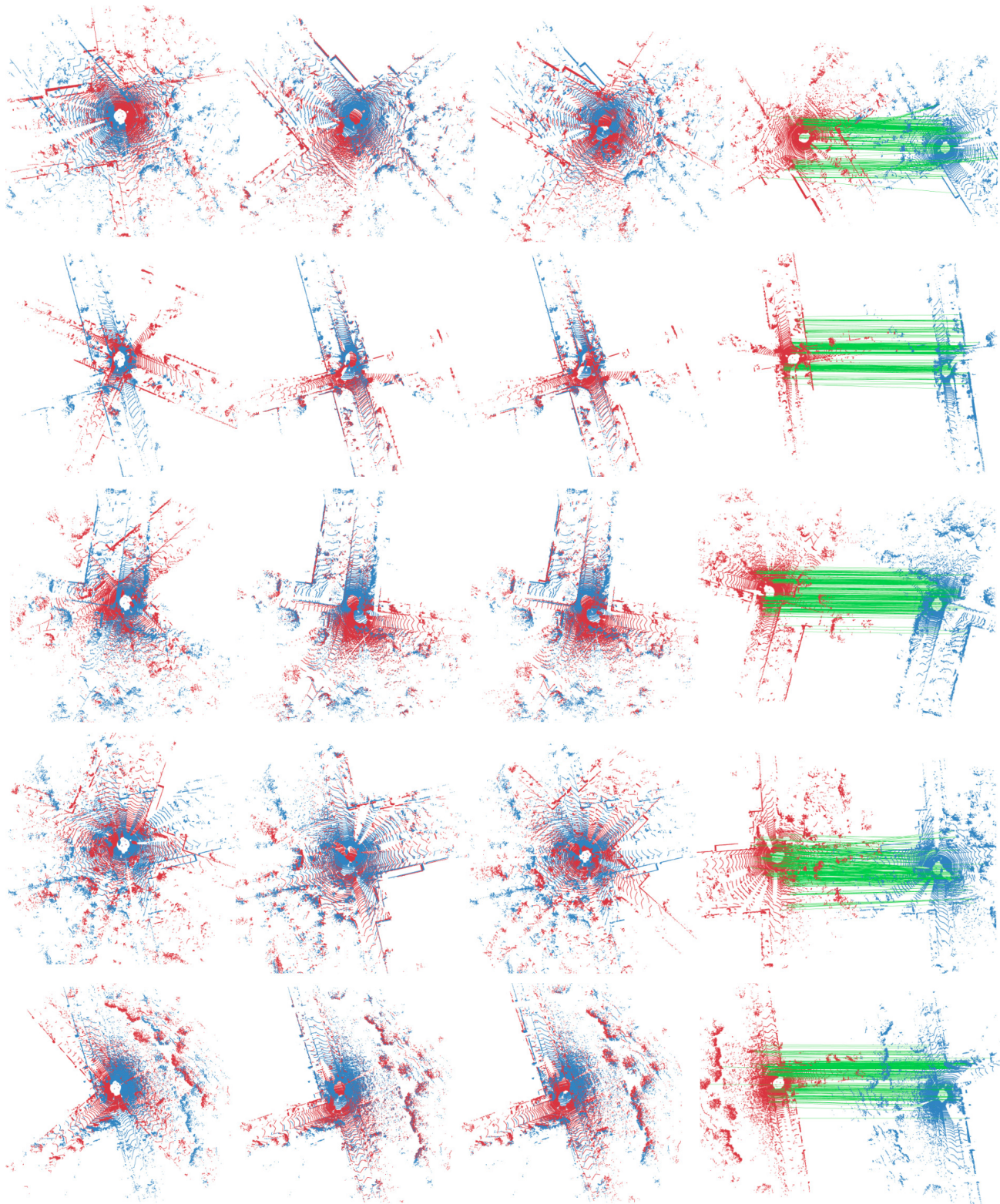


Figure 8. **Examples in bird-eye view of failed registrations on the KITTI odometry dataset [14].** From left to right: overlaid, non-registered input scans (blue and red colors); ground-truth registration; scans registered with our method; and top 256 pairs of matched points with highest confidence. *Note that we have used the full non-voxelized Lidar scans for a better visualisation. However, all registrations are done using 2048 points drawn at random after voxelization of the full point cloud.*

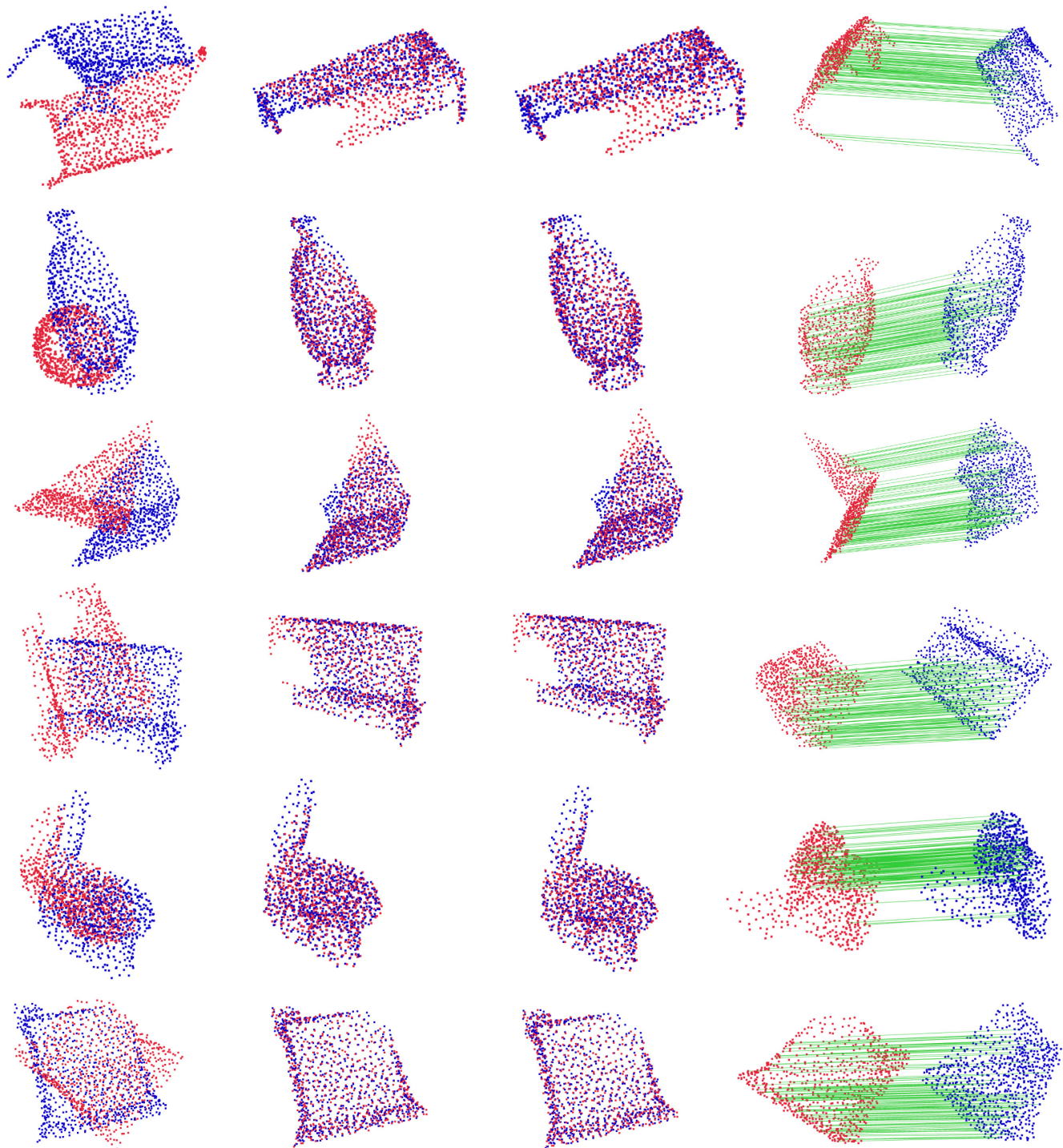


Figure 9. Examples of registration results on ModelNet40 on the split with unseen objects, with Gaussian noise [39]. From left to right: overlaid, non-registered input scans (blue and red colors); ground-truth registration; scans registered with our method; and top 128 pairs of matched points with highest confidence.