

Article

Have I seen this place before? A fast and robust loop detection and correction method for 3D Lidar SLAM

Michiel Vlaminck¹, Hiep Luong¹, and Wilfried Philips¹

- ¹ Department of Telecommunications and Information Processing, Ghent University, imec, Belgium
- * Correspondence: michiel.vlaminck@ugent.be; Tel.: +32-473413613
- + Sint-Pietersnieuwstraat 41, 9000 Gent, Belgium

Academic Editor: name Version November 21, 2018 submitted to Sensors

- Abstract: In this paper we present a complete loop detection and correction system developed for
- ² data originated from lidar scanners. Regarding detection, we propose to combine a global point
- ³ cloud matcher with a novel registration algorithm to determine loop candidates in a highly effective
- way. The registration method can deal with point clouds that are largely deviating in orientation
- ⁵ while improving the efficiency over existing techniques. In addition, we accelerated the computation
- of the global point cloud matcher by a factor 2 to 4, exploiting the GPU to its maximum. Experiments
- 7 demonstrated that our combined approach is more reliable to detect loops in lidar data compared
- * to other point cloud matchers as it leads to better precision-recall trade-offs: for nearly 100% recall,
- we gain up to 7% in precision. Finally, our loop correction algorithm is leading to an improvement
- ¹⁰ by a factor 2 on the average and median pose error while at the same time only needing a handful
- ¹¹ seconds to complete.
- ¹² **Keywords:** Loop detection; Lidar; Point Clouds.

13 1. Introduction

Mapping the environment using mobile mapping robots is a topic that by now has been studied 14 for almost two decades. Still, it can be considered as a highly active research area as the ultimate goal 15 of livelong mapping is far from reached. In the ideal case, livelong mapping comprises a solution in 16 which the map of the *world* is continuously updated at a pace equalling the one at which the world 17 is changing itself. During the past years many techniques have been introduced that can perform 18 incremental mapping using both regular cameras and depth sensing technologies based on either 19 structured light, ToF or pulsed lidar. As novel and more accurate sensors with increasing resolutions 20 are continuing to come to the market, the performance of these mapping solutions is still increasing. 21 However, even tough these sensing technologies are producing more accurate depth measurements, 22 they are still far from perfect and as a result, the proposed solutions suffer - and will continue to suffer 23 - from the drift problem. These drift errors could be corrected by incorporating sensing information taken from places that have been visited before. This requires both algorithms that can recognize 25 revisited areas as well as algorithms that can *close* the loop and propagate errors back in the *pose* 26 graph. Unfortunately, existing solutions for loop detection and closure for 3D data are computationally 27 demanding. In this work, we focus on speeding up both loop detection and loop correction in scanning 28 lidar data. We propose a technique that is able to automatically detect and correct *loops* in 3D lidar 29 data in a highly efficient way, thereby exploiting the power of modern GPU's. 30

31 2. Related work

Techniques to conduct loop detection in 3D data can roughly be categorized into three main 32 classes: local keypoint detection and matching in combination with a Bag-of-Words (BoW) approach, 33 global descriptor matching and a remainder category based on geometric primitives or whole objects. 34 The first class is generally detecting *salient* keypoints in a point cloud, computing *signatures* for 35 these keypoint positions, building a BoW and finally matching them in different scans [1]. Many 36 keypoint detectors have been proposed in literature, such as the intrinsic shape signatures (ISS) [2], 37 Harris 3D [3], Sift 3D [4], NARF [5], as well as many descriptors such as spin images [6], SHOT [7], 38 etc. However, despite this abundance of choice, the detection of distinctive keypoints with high 39 repeatability remains a challenging problem in 3D point cloud analysis. One way of dealing with this 40 lack of high repeatability is by using *global* descriptors which usually come in the form of histograms: 41 the fast point feature histogram (FPFH) [8] or Viewpoint Feature Histogram (VFH) [9] are a few examples. 42 Recently, He et. al. [10] presented a novel global 3D descriptor for loop detection, named multiview 43 2D projection (M2DP), that is very promising regarding both accuracy and efficiency. However, this global descriptor matching along with its local counterpart continues to suffer from respectively the 45 lack of descriptive power or the struggle with invariance. As a result, many loops are not detected or 46 too many false positives are present which on its turn imposes a restriction on fast and reliable loop 47 closure. More recently, researchers tend towards the application of convolutional neural networks 48 (CNNs) to learn both the feature descriptors as well as the metric for matching them in a unified 49 way [11–14]. A severe limitation of these methods on the other hand is that they need a tremendous amount of training data. Moreover, they don't generalize well when trained and applied on data 51 with varying topographies or acquired under different conditions. A model that was trained on data 52 originated from traffic environments, might perform poorly on indoor scenes and vice versa. Finally, 53 a third group of methods focus on place recognition based on either complete objects or planes. In 54 [15], Moral et. al. presented a place recognition algorithm based on plane-based maps. Unfortunately, their approach is only suitable for indoor man-made environments and is not scalable to outdoor 56 scenes. Dube et. al. [16] on the other hand proposed a method based on segments for which they 57 compute several descriptors that are integrated in a learning framework. Before feeding the segment 58 descriptors to the recognition model, they first subject them to a geometric verification test. As their 59 method is relying on all kinds of segments, it is more scalable than the work of Moral et. al.. However, when used in man-made environments such as buildings, some scans (e.g. parts of corridors) might 61 lack geometrical features and may jeopardize the loop detection. Besides the three aforementioned 62 categories, there are also 3D lidar SLAM systems that just use visual information to detect and close the 63 loop, such as the work of Zhu et. al. [17]. Regarding the actual loop correction, many solutions have 64 been presented in the literature as well [18–20]. Unfortunately, a severe drawback of these methods is their high complexity and computational burden, especially when the optimal solution is sought for. 66 Some algorithms therefore opt to conduct the loop closure in 2D, such as the work of Hess et. al. [21] in 67 which they propose a real-time loop closure algorithm for 2D lidar SLAM and which is part of Google's 68 cartographer. Another approach was taken in [22], where the authors propose a heuristic suitable for 69 large-scale 6D SLAM. Their idea is to conduct the optimization without any iteration between the 70 SLAM front- and back-end, yielding a highly efficient loop closing method. In this work we decided 71 to use the same strategy as we want to end up with an online - hence very fast - solution. In [23], 72 the authors propose to adopt a hierarchical approach instead by subdividing the 3D map into local 73 sub-maps. In order to incorporate corrections, the individual 3D scans in the local map are modelled 74 as a sub-graph and graph optimization is performed to account for drift and misalignments only at the 75 level of the local maps. 76

77 3. Contributions

The loop closure method developed in this work is part of an entire 3D mapping system that

⁷⁹ was presented in [24] and [25]. As in that work, we use lidar data originated from Velodyne scanners,

including the VLP-16, HDL-32E and HDL-64E. One of the main contributions of [24] was a scan
matching framework that aligns newly acquired point clouds with an online built 3D map. As stated
in the introduction this procedure is prone to error accumulation and for that reason we improve it by
actively detecting loops, i.e. locations that the robot or mobile platform is visiting more than once, in
order to propagate the error back in the SLAM pose graph. The main contributions of this work can be
summarized as follows:

 We accelerated the computation of a global 3D descriptor to detect strong loop candidates in lidar data. Compared to many other 3D descriptors, ours is not depending on the estimation of surface normals in the point cloud. The main motivation for this is the high difficulty of estimating these normals accurately given the sparse and inhomogeneous point density of lidar point clouds. Our 3D descriptor is thus leading to more robust loop detections. A compiled version of our algorithm has been made available to the community through a GitHub repository, allowing future use ¹.
 We propose a global registration technique based on 4-points congruent sets, inspired by the

work of Mellado *et. al.* [26]. We improved the efficiency significantly by omitting its randomized
 component which leads to faster execution times. Moreover, our improvements make the
 registration technique more robust for sparse and inhomogeneous 3D data.

3. We propose and evaluate a loop correction algorithm that omits the iteration between the SLAM

front- and back-end, leading to a very fast computations of the solution.

99 4. Approach

As mentioned in the previous section, our loop closure pipeline consists of two main steps. Prior 100 to these two steps, we additionally perform a quick selection of loop candidates by testing for each 101 newly computed position whether or not it is *close* to a place we already visited. By assuming that the 102 local registration is quite accurate, we set a threshold on the distance between two poses, equalling 10% of the total trajectory since the last loop closure. For instance, when we travelled for 100 meter 104 since the last loop closure, we consider a location to be matchable if its computed pose is within 10 105 meter. Once such a potential loop is detected, the next step deals with the matching of the 'start' and 106 'end' point cloud. As briefly mentioned in the previous section, this matching is done using a *global* 107 signature that we compute for each of the two point clouds. The signature is based on the projection of 108 the point cloud on several 2D planes, similar as the method described in [10]. If the matching residual 109 of this method is sufficiently low, we consider it as a strong loop candidate. To guarantee that we are 110 not dealing with a false positive, the next step seeks for the transformation that aligns the two point 111 clouds. In case the loop candidate is a true positive, we expect the overlap of the two point clouds 112 after registration to be very high. To this end, we adopt a *global* alignment technique, based on 4-points 113 congruent sets, to obtain a rough estimate of the transformation between the two ends of the loop. 114 The alignment of the two ends is essential to eventually close the loop, as it can happen that the same 115 position is revisited from an entirely different direction. Thus, the lidar scans can be acquired from 116 different viewpoints. The use of an ICP-based method is in this case not appropriate as it will converge 117 to a wrong local minimum. After the rough alignment, we still refine the transformation estimate 118 using a variant of the ICP algorithm. After this step we eventually do a quick geometric verification 119 check to see if the objects in the two (transformed) point clouds are relatively still located at the same 120 position. Figure 1 depicts a schematic overview to summarize our approach. 121

122 4.1. Multiview 2D-projection

Our global point cloud descriptor is inspired by the method of multiview 2D-projection (M2DP), first presented by He *et. al.* in [10]. The algorithm is summarized as follows. In order to achieve

¹ https://github.com/Shaws/m2dp-gpu



Figure 1. Our loop detection and correction pipeline, consisting of 5 main algorithms. The global descriptor MD2P detects loop candidates, after which the 4-PCS and ICP algorithm try to align both ends of the loop. Finally, the geometric verification step checks if the different segments in the scene are relatively still at the same position.

rotation invariance, we first compute the centroid of the point cloud - denoted by \mathcal{P} - and shift it 125 towards this centroid to achieve zero mean for the points. Second, we define a reference frame by 126 estimating two dominant directions of the point cloud using PCA. The two principal components are 127 then set as the x- and y-axis of the descriptor reference frame. The third step consists of generating 128 several 2D signatures by defining different planes on which we project the 3D data. Each plane can 129 be represented using the *azimuth* angle θ and *elevation* angle ϕ . Thus, each pair of parameters $[\theta, \phi]$ 130 leads to a unique plane X_i with normal vector $\mathbf{n}_i = [\cos\theta\cos\phi, \cos\theta\sin\phi, \sin\theta]^{\dagger}$. The projection of a 131 point $\mathbf{p}_i \in \mathcal{P}$ on X_j is then given by $\mathbf{p}_i^j = \mathbf{p}_i - \frac{\mathbf{p}_i^\top \mathbf{n}_j}{||\mathbf{n}_j||_2^2} \mathbf{n}_j$. To describe the structure of the points on X_j , 132 the 2D plane is further divided into bins as follows (cfr. Figure 2). Starting from the projected centroid 133 on X_i , *l* concentric circles are generated, with radii $[r, 2^2r, l^2r]$ and the maximum radius is set to the 134 distance of the centroid and the furthest point. Each ring is divided in t bins, hence defining $l \times t$ 135 different bins. For every bin, the number of projected points lying in it are counted, generating a $lt \times 1$ 136 signature vector v_i describing the points projected on plane X_i . The main benefit of this projection is 137 that it is not relying on the estimation of surface normals for the points. This latter procedure is usually 138 time consuming and often times inaccurate for lidar data given their sparse and inhomogeneous point 139 density. The signature is computed for p different azimuth angles and q different elevation angles 140 where the stride on azimuth is $\frac{\pi}{p}$ and the one on elevation $\frac{\pi}{2q}$. Hence, there are pq different planes, 141 leading to a signature matrix A of size $pq \times lt$, for which each row corresponds with a single signature 142 vector v_i . Finally, a SVD decomposition is run on the matrix A and the first left and right singular 143 vectors are concatenated to form the final descriptor. Many parts of this algorithm lend themselves 144 to be computed in parallel, allowing us to exploit the multi-core nature of modern GPU's. To ease 145 the implementation, we made use of Quasar, a language and computing platform facilitating GPU 146 programming that was presented in [27]. Specifically, three major parts were accelerated, the first one 147 being the determination of the two dominant directions of the point cloud computed using PCA. Our 148 PCA implementation uses a parallelized version of the SVD algorithm. The second part deals with the 149 projection of the points on the different planes X_i , $j = 1 \dots pq$, all of which can be computed in parallel. 150 This leads to a speed up for this part of $pq \times N$, N being the number of points in \mathcal{P} , compared to a 151 serial version of the algorithm. Finally, once the projections are computed for all planes and all points, 152 the numbers of points belonging to each bin can also be determined in parallel. 153



Figure 2. The generation of a 2D signature by projecting the 3D data onto a plane. The plane is divided into bins as follows: starting from the projected centroid, *l* concentric circles are generated and the maximum radius is set to the distance between the centroid and the furthest point. Each ring is subsequently subdivided into *t* bins, generating a $lt \times 1$ signature vector. Finally, the number of points lying in each bin are counted.

154 4.2. 4-points congruent sets

The algorithm described in the previous section provides us some loop candidates. However, 155 its outcome is insufficient to be able to close the loop. As the descriptor is rotation invariant, we 156 need to find out the transformation between the two point clouds in order to *discard* it from the error 157 back propagation. To this end, we need a registration algorithm that can deal with large variations in orientation. Figure 3 depicts a bird's-eye view of the point clouds at both ends of the loop and 159 demonstrates why the registration is necessary. The black and gray point clouds are respectively the 160 start and end point of a loop. The green point cloud is the transformed version of the gray one. If 161 after registration the overlap between the two point clouds (green and black) is sufficiently high, we 162 accept it is a true loop. Another key thing to keep in mind is that it can happen that a part of the scene 163 has been changed in the mean time. Think about parked cars along the road that disappear or other 164 ones that have taken their place. The use of local feature descriptors in a BoW approach would most 165 likely fail in these cases. For that reason, our global registration technique takes into account the main 166 geometry of the scene instead of relying on keypoint positions. The core of our registration technique 167 is based on 4-points congruent sets (4-PCS), an idea initially proposed by Aiger et. al. [28] and later 168 improved by Mellado et. al. [26]. Our 4-PCS method is thus a global point cloud registration technique 169 that does not rely on the extraction of features. Instead it is matching congruent sets in both point 170 clouds thereby adopting a generate-and-test paradigm, known from RANSAC-based solutions. In its 171 most simple form, RANSAC randomly selects three points from both the source point cloud \mathcal{P} and the 172 target point cloud ${\cal Q}$ and subsequently computes the rigid transformation using these points. Next, 173 it tries to *verify* this transformation by determining how many points from \mathcal{P} are within a δ -distance 174 from points in Q after applying the transformation. If this count - usually referred to as the size of the 175 consensus set - is sufficiently high, the transformation is accepted as a good solution. Otherwise, the 176 process is repeated by randomly selecting another triplet of points. The transformation with the largest 177 consensus set is finally accepted as the best fit. The 4-PCS method builds on this randomized alignment 178 approach, but instead of exhaustively testing all the triplets from Q, it introduces the concept of *planar* 179



Figure 3. Bird's-eye view of two locations originated from the KITTI 00 sequence. The black and gray scans have been selected by the M2DP algorithm as the start and the end point of a potential loop. The two point clouds are however slightly rotated. Our registration technique aligns both (the green one is the transformed version of the gray point cloud) and determines their overlap. When the overlap is sufficiently high, the loop candidate is eventually selected as a true loop.

congruent sets to select only a small subset of bases from Q that can potentially match a given base 180 from \mathcal{P} . The first step in the 4-PCS method thus consists of selecting a 4-point coplanar base B from 181 the source point cloud \mathcal{P} . Next, from the target point cloud \mathcal{Q} , all 4-point sets $\{U_1, \ldots, U_N\} = \mathbf{U}$ that 182 are approximately congruent to B are determined. Third, for all sets U_i , the rigid transformation T_i 183 that aligns B and U_i is computed and verified according to the *largest common point set* (LCP) criterion. This latter criterion denotes the set of points $S_i \in \mathcal{P}$ that are within δ -distance from a point in \mathcal{Q} after 185 applying the transformation. Finally, the best transformation T_{opt} , i.e. the one leading to the set S_k 186 with the highest cardinality is kept. In summary, the aforementioned algorithm consists of four major 187 steps: 1) selecting a coplanar base in one point cloud, 2) find the (approximate) congruent sets in the 188 second point cloud, 3) compute the rigid transformations and 4) test the rigid transformations and select the best one. 190

4.2.1. Selecting a coplanar base

One of the main limitations of the original method by Aiger *et. al* [28] is its random search for coplanar points which is leading to a lot of unnecessary computational burden. Instead, we propose to 193 cluster the 3D points and subject them to a plane fitting process. As the point clouds are generated by 194 a scanning lidar device with 16 to 64 colinear lasers, we can project the 3D laser points onto a regular 195 2D grid, as described in [24]. Doing so, we can exploit the known adjacency of the points to quickly 196 perform clustering. More specifically, we adopt a region growing algorithm using two comparator functions to determine whether or not two neighbouring points are belonging to the same cluster. 198 The first one is the Euclidean 3D distance between the two points, the second one the deviation of 199 their surface normal. After applying this region growing process we obtain a set of clusters, which we 200 subsequently feed to a plane fitting algorithm. Once we have eventually found some clusters to be 201 part of a plane, we can extract coplanar bases very easily as any four points lying in the same plane are 202 by definition coplanar. Following this procedure, we can omit the randomized base selection process 203 of the original 4-PCS method. Obviously, it is still beneficial to pick wide bases (by selecting points 204 that are located far from each other) as they are in general leading to more stable alignments. To this 205 end, we prioritize points lying at the boundaries of the planar cluster to serve as a base. Of course, the 206 base should still lie in the overlap region between the two point clouds in order not to miss the desired 207 solution. As we are considering point clouds that are captured at more or less the same position (but at 208 different moments in time) we assume that the overlap will be quite large. Only in the case that a large object close to the scanner is causing a huge occlusion in one of the point clouds, this assumption might 210 be violated. Therefore, we propose to compute for each planar region its convex hull and to select 211 a coplanar base by picking four points that are close to this convex hull. Figure 4 depicts two point 212



Figure 4. The process of selecting a coplanar base. First, several objects in the scene are clustered after which they are subjected to a plane fitting algorithm. In case a few planes have been detected, their convex hulls are computed. A coplanar base is then selected by picking four points close to the convex hull.

clouds from sequence '05' of the KITTI benchmark together with the convex hulls of the estimatedplanar regions.

²¹⁵ 4.2.2. Finding congruent sets

Once a coplanar base is selected, the next step consists of finding 4-points sets in the other point cloud that are congruent to this base. This matching step is based on a specific property of affine invariants of 4-points congruent sets. In a nutshell, given 4 points, we can compute two independent ratios between the line segments they are defining that are preserved under affine transformations. Given a set of coplanar points $\mathbf{B} = \{\mathbf{a}, \mathbf{b}, \mathbf{c}, \mathbf{d}\}$ from point cloud \mathcal{P} that are not all collinear. Let \mathbf{ab} and \mathbf{cd} be the two lines that intersect at an *intermediate* point *e*. The two ratios

$$r_1 = \frac{||\mathbf{a} - \mathbf{e}||}{||\mathbf{a} - \mathbf{b}||}, \qquad r_2 = \frac{||\mathbf{c} - \mathbf{e}||}{||\mathbf{c} - \mathbf{d}||} \tag{1}$$

are invariant under affine transformation, and uniquely define 4-points up to affine transformations. Now, for each point $\mathbf{q}_1, \mathbf{q}_2 \in \mathcal{Q}$, we can compute two *intermediate* points:

$$\mathbf{e}_1 = \mathbf{p}_1 + r_1(\mathbf{p}_2 - \mathbf{p}_1), \qquad \mathbf{e}_2 = \mathbf{p}_1 + r_2(\mathbf{p}_2 - \mathbf{p}_1).$$
 (2)

Any two pairs whose intermediate points \mathbf{e}_1 and \mathbf{e}_2 are coincident, potentially correspond to a 4-points 216 set that is an affine transformed copy of **B**. Of course, as these sets of 4-points are the affine invariants 217 of the base **B**, it is a superset of the 4-points set that are a rigid transformation of the base. For that 218 reason, we also check the angle between the two line segments to determine if the 4-points set is a 219 rigid transformation of the base **B**. Naturally, the intermediate points e_1 and e_2 will never exactly be 220 coincident due to noise and other inaccuracies. Instead, they will end up on being *nearby* points. For 221 that reason, we set up a k-d tree search data structure that allows for fast spatial queries. We then 222 accept the set as being congruent to the base **B** in case the distance of the two intermediate points \mathbf{e}_1 223 and \mathbf{e}_2 are within δ -distance from each other. Another limitation of the original method of [26] is that 224 it computes and tests for all possible combinations of points their intermediate points \mathbf{e}_1 and \mathbf{e}_2 . This 225 is leading to a tremendous amount of unnecessary computations. Instead, we propose to only process 226 the points of clusters that are lying in a physical plane, i.e. a plane that is present in the scene. Only 227 these points qualify to be matchable with a given coplanar base, as the bases themselves were picked 228 on the detected planar regions. 229

230 4.2.3. Test rigid transformation

The final step in the 4-PCS method is to test the rigid transformation computed using the base and its congruent sets. One way of verifying the transformation is by using the *largest common point set* (LCP) criterion. This criterion states that one should count the number of points from the *source*

point cloud that are within a δ -distance from any point from the *target* point cloud after alignment. The 234 transformation that yields the largest LCP is considered as the true transformation. We emphasize that 235 we only compute the transformation and LCP criterion for the congruent sets that we have selected in the previous step. This group of congruent sets is thus a lot smaller than the original method proposed 237 by [26] et. al.. Our algorithm can thus be summarized as follows. First, select a few strong coplanar 238 bases from \mathcal{P} based on the estimated planes in the scene. Second, given a selected coplanar base from 239 \mathcal{P} , determine only a few strong 4-points sets in \mathcal{Q} that are approximately congruent. For all these 240 selected 4-points sets, compute the transformation that aligns the two point clouds and eventually pick the transformation with the largest common point set. In other words, pick the transformation for 242 which the most points of \mathcal{P} are within δ -distance from a point in \mathcal{Q} . 243

4.3. Verification through ICP and geometrical consistency

The 4-PCS method yields a rough transformation from source to target point cloud but it will not 245 perfectly align them. We therefore refine the transformation by using an ICP-based algorithm that 246 was described in [24]. This ICP algorithm offers an additional verification to conclude that the two 247 point clouds yield a true loop. If the residual of ICP - defined as the average distance between all 248 corresponding points - is too large we still discard the loop candidate. Finally, the object clustering mentioned in previous section also provides a means of verification as we can check if all the object 250 clusters are relatively still at the same position. To that end, we compute for each cluster its centroid 251 and compute a bipartite matching using the Hungarian algorithm. In case these two verification steps 252 are positive we can eventually proceed to the actual loop correction. 253

254 4.4. Loop correction

As stated in the introduction our goal is to implement loop closure as an online process. Therefore, we adopt a heuristic approach that is extremely fast, though leading to a sub-optimal solution. The idea is to avoid the iteration between the SLAM front- and back-end. The front-end is referring to the scan matching process whereas the back-end deals with the global consistency of the 3D map and hence the correction of the loops. In that iteration, the outcome of the SLAM back-end, being the pose error, is given to the front-end to re-investigate its outcome thereby taking into account all the known relations between neighbouring poses and matched correspondences. After this, the outcome is feeded back to the back-end to check if this re-investigation has lead to a better result. This process is repeated for all poses and all correspondences, up untill the optimal solution is found. It goes without saying that all this is inherently time consuming and therefore we propose to bypass this iterative behaviour. Instead, we pass the *local* information from the scan matching just once to the back-end, but not the other way round. In a nutshell, we investigate how much each pose is contributing to the final accumulated error, thereby considering the residual score of the scan matching process. The actual correction is then done as follows. Consider the mapping platform travelling along the trajectory $V_0, \ldots, V_1, \ldots, V_n$ where at each pose V_i the lidar scanner is capturing a point cloud \mathcal{P}_i during a full rotation of its head, also referred to as a sweep. We assume that the loop detection method provided us with two point clouds, resp. the *start* and *end* of the loop and let us denote their poses as respectively V_s and V_e . Next, the difference in pose, i.e. the loop transform, is given by $\Delta = (R_{s,e}V_e)^{-1}V_s$. Note that the rotation $R_{s,e}$ denotes the one we have computed using our registration algorithm. It should be discarded from the loop correction process as it does not yield an error. This loop transform Δ is considered as an *error* as both poses $R_{s,e}V_e$ and V_s should be equal. It should thus be projected back in the pose graph. As mentioned before, we use the residual of the scan matching process, i.e. final cost after transformation, to assign a weight $c_{i,j}$ to each edge in the pose graph. We assign a higher weight for those transformations that yield a high residual in the scan matching step. The idea is that a high residual indicates that two consecutive point clouds were potentially inaccurately aligned. In addition, we assume that the scan matching process will have already been converging in the right direction.

Table 1. Computation times of our GPU implementation of the M2DP descriptor compared to the CPU version implemented in Matlab by [10]. The speed-up factor for data originated from the HDL-64E is almost x4 whereas for the VLP-16 and HDL-32E the speed-up factor is rather x2.5.

Sensor	avg. $ \mathcal{P} $	Matlab ([10]) in ms	GPU (ours) in ms
VLP-16	25446	250	110
HDL-32E	51687	273	121
HDL-64E	62594	476	127

Next, we define the *distance* between two poses V_k and V_l as $d(V_k, V_l) = \sum_{i,j} c_{i,j}$. Herein, $\{i, j\}$ denotes the set of all edges in the path from V_k to V_l . Finally, we define a weight

$$w_i = \frac{d(V_s, V_i)}{d(V_s, V_e)} \tag{3}$$

for each pose in the graph that specifies the fraction of the matrix Δ by which the pose has to be transformed. The poses V_k are than updated replacing \mathbf{t}_k by $\mathbf{t}_k w_k \Delta$ and R_k by $slerp(\mathbf{R}_k, w_k \Delta)$, slerpdenoting the spherical linear interpolation function as described in [29].

258 5. Evaluation

The evaluation covers three main parts. First, we conduct an analysis to qualify the speed-up of our GPU-accelerated descriptor computation. Second, we compare our loop detection accuracy to other state-of-the-art methods. Finally, the quality and the speed of the loop closure module are analysed.

263 5.1. Speed analysis global descriptor computation

In order to compare our GPU implementation of the M2DP descriptor against the CPU version 264 implemented in Matlab by [10], we ran several experiments using data originated from several 265 Velodyne scanners. To indicate the performance on point clouds of different sizes, we used data from 266 both the VLP-16, HDL-32E and HDL-64E containing resp. 16, 32 and 64 lasers. More specifically, we 26 used several sequences from the KITTI benchmark [30] which are all captured by a HDL-64E in urban 268 environments. We extended this dataset with own recorded sequences in both indoor and outdoor 269 scenes using the HDL-32E and VLP-16. The experiments were conducted on a computer with an Intel 270 Core i7-7820X @ 3.60Ghz, 128GB RAM and an nVidia GeForce GTX 1080ti inside. The results are 271 summarized in Table 1. As can be seen, our GPU implementation scales well for larger point clouds. 272 The overhead of copying data to the GPU memory is relatively lower for larger point clouds, hence 273 yielding a larger speed-up. For point clouds originated from the HDL-64E scanner our implementation 274 only takes 127 ms on average compared to 476 ms for the Matlab implementation of [10] which is a 275 speed-up factor of nearly 4. For smaller point clouds - captured with the VLP-16 or HDL-32 - we notice 276 a speed-up factor of almost 2.5. 277

278 5.2. Loop detection accuracy

In order to evaluate the accuracy of our *final* loop detector against the state-of-the-art, we 279 conducted several experiments on the Kitti benchmark [30]. More specifically, we used the sequences 280 '00', '05' as these contain the most 'revisited' locations. To generate ground truth we used the known 281 282 trajectories and considered a loop to be present when the distance between two poses is less than 1 meter. Thereby, we used a threshold of 100 poses to avoid that two subsequent poses would wrongly 283 be classified as a loop. In Figure 5 the 'ground truth' revisited locations in the two trajectories are 284 depicted as green dots. For some sequences the same road has been taken multiple times, hence the 285 whole part of the road is considered as a loop. However, as can be seen, sometimes there are green 286



Figure 5. The ground truth loops in the sequences '00' and '05' of the Kitti benchmark [30]. We consider a loop to be present when the distance between two poses is less than 1 meter.

dots missing along a road that has been taken multiple times, meaning that the difference in pose is 287 larger than 1 meter. This could be due to the fact that some roads consist of multiple lanes and that a 288 different lane was taken. We used four different descriptors to compare our results with. The first one 289 is the original M2DP method described in [10]. The other descriptors are the *ensemble of shape functions* 290 (ESF) [31], spin images [32] and SHOT [7]. The latter two are local descriptors so in order to use them 29: as a global descriptor we computed the centroid of the point cloud and estimated the spin image and 292 SHOT descriptor related to this centroid. Thereby we use the maximum distance from the centroid 293 to any other point in the cloud as the radius to compute the descriptor. Furthermore, both the spin 294 images and the SHOT descriptors are based on normal vectors. To compute these, we use a radius that 295 equals five times the average distance of a point to its closest neighbour. Regarding spin images, the 296 other parameters that need to be set are 1) the number of bins along one dimension, 2) the minimal 297 allowed cosine of the angle between the normals of the input cloud and search surface (for the point to 298 be retained in the support) and 3) the minimal number of points in the support to correctly estimate 200 the spin image. We have set these parameters to respectively 8, 0.5 and 16. The SHOT descriptor does 300 not have any other parameters to be set and for the ESF we again used the maximum distance of the 301 centroid to any other point in the cloud. Finally, the M2DP method needs the number of bins for each 302 plane (or expressed as the number of circles l and number of bins in one ring t) and the number of 303 planes to use (or expressed as the azimuth *p* and elevation *q*). For our experiments we have set these 304 values to l = 8, t = 16, p = 4 and q = 16 for all tests. The ROC curves for the different loop detectors 305 are depicted in figure 6. We clearly see that our method along with the methods M2DP and SHOT 306 are leading to the best detections. The performance of the ESF descriptor and spin images turns out to be insufficient to reliably recognize revisited areas. To obtain a recall of at least 90%, the precision 308 drops to respectively 45% and 20% for the KITTI sequence '00', which is unacceptable in an operational 309 system. For the KITTI sequence '05' the precisions corresponding to a recall of 90% are even worse, 310 resp. 30% and 15%. On the contrary, the M2DP and SHOT descriptor are leading to a precision of 70% 311 on the '00'-sequence and higher than 95% on the '05'-sequence for a 90% recall. In table 2, the exact 312 precision is listed that corresponds to a recall of about 99.9%. We observe that our combined method 313 further improves the performance of the M2DP detector. For the KITTI sequence '05', to obtain 99.9% 314 recall, we reach a precision of 90.4%, an improvement of more than 7.1%. For the KITTI sequence 315 '00', we obtain a smaller improvement of 0.8% reaching a precision of 60.5% compared to 59.7 for the 316



Figure 6. ROC curves for the several loop detection methods on the KITTI sequences '00' and '05'. As can be seen, our method performs better than the original M2DP method, the second best performer in our experiments. The SHOT detector also produces acceptable results, whereas the spin and ESF detectors are too unreliable to be used in practise.

Table 2. Precision at 99.9% recall of the different loop detection methods for the KITTI dataset. Clearly, only the detectors SHOT and M2DP are producing 'acceptable' results. Our combined algorithm improves the original M2DP method by 0.8% and 7% for KITTI sequence '00' and '05' leading to a precision of respectively 60.5% and 90.4%.

Sequence	Spin Image	SHOT	ESF	M2DP	ours
KITTI00	0.025	0.575	0.143	0.597	0.605
KITTI05	< 0.01	0.632	0.037	0.833	0.904

M2DP detector. The lower performance on the KITTI '00' dataset is probably due to the inaccuracies 317 in the ground truth. It goes without saying that a higher threshold on the distance for a pose to be 318 considered as a ground truth loop affects the results tremendously. Thus, the value of this experiment 319 is in the comparison between the different detectors rather than in the absolute numbers on the actual accuracy. Besides experiments on the Kitti benchmark, we acquired a lidar sequence ourselves in the 321 city of Ghent (Belgium) thereby mounting the Velodyne VLP-16 lidar scanner on top of a car. While 322 acquiring the lidar data we used a Garmin GPS to generate ground truth. The trajectory is shown in 323 Figure 7. As can be seen, a part of the trajectory was taken twice, making all the poses along this part 324 act as loops. This time we used 3 meters as a threshold for two poses to be considered as a loop as this threshold was leading to more coherent loops along roads that were taken twice. As the sequence 326 was recorded in the historical city centre of Ghent, the GPS signal was often times inaccurate leading 327 to a noisy trajectory. To deal with these anomalies, we used the Google API to *clean-up* the trajectory 328 by computing the most likely roads that were taken. This eventually lead to resp. 843 loops on a 329 trajectory of approx. 15.7km travelled in 47 minutes. The number of lidar point clouds is 31745. For 330 this experiment we used the exact same parameters as for the KITTI sequences. The ROC curve is 331 depicted in Figure 7. In this experiment, our loop detector is clearly outperforming the other point 332 cloud matchers. For a recall of 93% we still obtain a precision of 100%, which is an improvement of 333 5.2% compared to the SHOT descriptor. Only the latter along with the M2DP descriptor are generating 334 acceptable results. The spin images are performing better than for the KITTI sequences, but the overall 335 accuracy is still too low to be usable in an operational system. Finally, for the ESF descriptor it is very



Figure 7. Left: the ground truth and estimated loops of the sequences recorded in the city center of Ghent, Belgium. We consider a loop to be present when the distance between two poses is less than 3 meter. Right: the ROC curve for the sequence recorded in the city center of Ghent, Belgium. Our loop detector is clearly outperforming the other methods. For a recall of 93% we still obtain a precision of 100%, which is a 5.2% gain compared to the SHOT descriptor. This latter along with the M2DP descriptor are the only two other point cloud matchers that are producing acceptable results. The loop detection quality of the spin images descriptor and ESF are too limited.

difficult to achieve any acceptable precision, even for recall values. Due to inaccuracies of the GPS
data, and hence the ground truth, it was not possible to reach 100% recall for any of the methods, as
was the case for the KITTI.

5.3. Quality and speed of the loop correction

To assess the quality of our loop correction algorithm, we conducted an experiment using data that we captured in the Belgian city Hasselt. We mounted a Velodyne HDL-32e on a mobile mapping van together with a high-precision POS-LV inertial positioning system to acquire accurate *ground truth*. In order to obtain the *initial* trajectory we used our 3D reconstruction system described in [24]. In Figure 8, two images are depicted showing the result of this reconstruction process. In the left image, the resulting point cloud is shown before loop closure, whereas the right picture depicts the point cloud after loop closure. As can be seen in the left point cloud, the end of the loop is not connected



Figure 8. The two resulting reconstructions for the 'Hasselt' dataset, before loop closure (left) and after loop closure (right). After loop correction, the two ends of the loops are attached and the error is propagated back in the pose graph.



Figure 9. The two resulting trajectories for the 'Hasselt' dataset corresponding to the reconstructions in Figure 8, before loop closure (left) and after loop closure (right). The estimated trajectory and the ground truth trajectory are almost entirely overlapping after loop closure. The pose error before and after loop correction is respectively 5.56 and 3.13 meter, hence proving the effectiveness of the loop correction.

with the start of the loop. After correction the two ends are connected and the accumulated error 348 is propagated back in the pose graph. In order to evaluate this quantitatively, we used the POS-LV 349 positioning system as ground truth. Figure 9 depicts both the ground truth and the *final* trajectory after 350 performing loop closure. As we do not have a complete ground truth 3D model, we measured the 351 quality by means of comparing all the poses in the pose graph. To this end we computed the average distance from all poses in the estimated trajectory with its *closest* pose in the ground truth. The average 353 distance between all poses before and after loop closure turned out to be respectively 5.65 meter and 354 3.13. Furthermore, the median pose error before and after loop closure are respectively 3.98 meter and 355 2.13. Hence, the loop closure process reduces the total error almost by a factor 2. The total time to 356 correct the loop on an Intel Core i7-4712HQ CPU @ 2.30 GHz, was 11 milliseconds. This is extremely 357 fast thanks to the omission of the iteration between the SLAM front- and back-end. During the SLAM 358 front-end a residual score was computed that is later on used in the loop correction phase. The loop 359 correction itself only involves one single manipulation of the poses. By means of a second example, 360 we also closed the loop of sequence '09' of the Kitti benchmark. The result is shown in Figure 10. The 361 trajectory after loop closure (in green) and the 'ground truth' trajectory (in red) are almost entirely 362 overlapping. On the right image, the point cloud at the start and end-point of the loop are shown. 363 The poses are overlapping and little to no artefacts can be seen in the point cloud. For this dataset we 364 also computed the average pose error before and after loop correction, which resulted in a value of 365 respectively 9.89 meter and 4.80 meter. The median value before and after loop correction is 7.85 meter 366 and 3.53 meter respectively. In summary, we can conclude that we reduce the error with a factor 2 after 367 loop closure. The results of both experiments are summarized in Table 3. 368

369 6. Conclusion

In this paper a full loop detection and correction method was presented. The main contributions of this work were twofold. First, we accelerated the loop detection process by developing a GPU-accelerated version of a global feature descriptor for point clouds. Second, we presented a novel registration technique to align two point clouds with a large deviation in orientation. We improve over existing techniques regarding both robustness, accuracy and speed. We explained



Figure 10. The resulting trajectory (left) of sequence '09' of the Kitti benchmark, before loop closure (blue) and after loop closure (green). The estimated and the ground truth trajectory are almost entirely overlapping after loop correction. Right: the corresponding reconstruction at the start and end-point of the loop after correction. Little to no artefacts can be seen in the reconstruction. The pose error before and after loop correction is respectively 9.89 and 4.80 meter, hence proving its effectiveness.

Table 3. The average and median pose error before and after loop closure for the Hasselt and Kitti '09' sequence. Our correction algorithm reduces the error by a factor of approximately **x2** after closing the loop.

Sequence	total length (m)	Average error (m)			Median error (m)		
		before LC	after LC	gain	before LC	after LC	gain
HASSELT	715	5.65	3.13	x1.81	3.98	2.13	x1.87
KITTI09	1705	9.89	4.80	x2.06	7.85	3.53	x2.22

how this method can be used to align two point clouds that serve as the two ends of a loop in lidar
data. Furthermore, we showed the effectiveness of incorporating this registration technique in the
verification process of potential loop candidates. Experiments have demonstrated that we can gain up
to 7% precision for the same recall value of 99.9%. Experiments also showed that our method works
for data acquired with different Velodyne scanners, i.e. the one containing 64 lasers as well as the one
containing only 32 or 16 lasers. In addition, we showed that our global feature descriptor is a factor
2.5 to 4 times faster than the original version, depending on the number of points in the point cloud.
Finally, we showed that our loop correction method reduces the average and mean pose error - defined

as the distance of a pose with its closest neighbour in the ground truth trajectory - by a factor 2.

384 References

- Steder, B.; Ruhnke, M.; Grzonka, S.; Burgard, W. Place recognition in 3D scans using a combination of bag
 of words and point feature based relative pose estimation. IROS. IEEE, 2011, pp. 1249–1255.
- Zhong, Y. Intrinsic shape signatures: A shape descriptor for 3D object recognition. 2009 IEEE 12th
 International Conference on Computer Vision Workshops, ICCV Workshops, 2009, pp. 689–696.
- Sipiran, I.; Bustos, B. A Robust 3D Interest Points Detector Based on Harris Operator. Eurographics
 Workshop on 3D Object Retrieval; Daoudi, M.; Schreck, T., Eds. The Eurographics Association, 2010.
- Scovanner, P.; Ali, S.; Shah, M. A 3-dimensional Sift Descriptor and Its Application to Action Recognition.
 Proceedings of the 15th International Conference on Multimedia; ACM: New York, NY, USA, 2007;
 MULTIMEDIA '07, pp. 357–360.
- Steder, B.; Rusu, R.B.; Konolige, K.; Burgard, W. NARF: 3D Range Image Features for Object Recognition.
 Workshop on Defining and Solving Realistic Perception Problems in Personal Robotics at the IEEE/RSJ Int.
 Conf. on Intelligent Robots and Systems (IROS); , 2010.
- Johnson, A. Spin-Images: A Representation for 3-D Surface Matching. PhD thesis, Robotics Institute,
 Carnegie Mellon University, Pittsburgh, PA, 1997.
- Salti, S.; Tombari, F.; di Stefano, L. SHOT: Unique signatures of histograms for surface and texture
 description. *Computer Vision and Image Understanding* 2014, 125, 251–264.
- Rusu, R.B.; Blodow, N.; Beetz, M. Fast Point Feature Histograms (FPFH) for 3D Registration. Proceedings
 of the 2009 IEEE International Conference on Robotics and Automation; IEEE Press: Piscataway, NJ, USA,
 2009; ICRA'09, pp. 1848–1853.
- Rusu, R.B.; Bradski, G.R.; Thibaux, R.; Hsu, J.M. Fast 3D recognition and pose using the Viewpoint Feature
 Histogram. IROS. IEEE, 2010, pp. 2155–2162.
- He, L.; Wang, X.; Zhang, H. M2DP: A novel 3D point cloud descriptor and its application in loop closure
 detection. IROS. IEEE, 2016, pp. 231–237.
- Dewan, A.; Caselitz, T.; Burgard, W. Learning a Local Feature Descriptor for 3D LiDAR Scans. Proc. of the
 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS); , 2018.
- I2. Zeng, A.; Song, S.; Nießner, M.; Fisher, M.; Xiao, J.; Funkhouser, T. 3DMatch: Learning Local Geometric
 Descriptors from RGB-D Reconstructions. CVPR, 2017.
- Yin, H.; Ding, X.; Tang, L.; Wang, Y.; Xiong, R. Efficient 3D LIDAR based loop closing using deep neural network 2017. pp. 481–486.
- Yin, H.; Wang, Y.; Tang, L.; Ding, X.; Xiong, R. LocNet: Global localization in 3D point clouds for mobile
 robots. *CoRR* 2017, *abs*/1712.02165, [1712.02165].
- Fernández-Moral, E.; Mayol-Cuevas, W.W.; Arévalo, V.; Jiménez, J.G. Fast place recognition with
 plane-based maps. ICRA. IEEE, 2013, pp. 2719–2724.
- ⁴¹⁸ 16. Dubé, R.; Dugas, D.; Stumm, E.; Nieto, J.I.; Siegwart, R.; Cadena, C. SegMatch: Segment based loop-closure
 ⁴¹⁹ for 3D point clouds. *CoRR* **2016**, *abs*/1609.07720.
- 420 17. Zhu, Z.; Yang, S.; Dai, H.; Li, F. Loop Detection and Correction of 3D Laser-Based SLAM with Visual
- Information. Proceedings of the 31st International Conference on Computer Animation and Social Agents;
 ACM: New York, NY, USA, 2018; CASA 2018, pp. 53–58.
- 18. Kümmerle, R.; Grisetti, G.; Strasdat, H.; Konolige, K.; Burgard, W. G2o: A general framework for graph
 optimization. ICRA. IEEE, 2011, pp. 3607–3613.

- Grisetti, G.; Stachniss, C.; Burgard, W. Non-linear Constraint Network Optimization for Efficient Map
 Learning. *IEEE Transactions on Intelligent Transportation Systems* 2009, 10, 428–439.
- Lu, F.; Milios, E. Globally Consistent Range Scan Alignment for Environment Mapping. *Auton. Robots* **1997**, 4, 333–349.
- Hess, W.; Kohler, D.; Rapp, H.; Andor, D. Real-Time Loop Closure in 2D LIDAR SLAM. 2016 IEEE
 International Conference on Robotics and Automation (ICRA), 2016, pp. 1271–1278.
- 431 22. Sprickerhof, J.; Nüchter, P.A.; Lingemann, K.; Hertzberg, P.J. A Heuristic Loop Closing Technique for
 432 Large-Scale 6D SLAM. *Automatika* 2011, 52, 199–222, [https://doi.org/10.1080/00051144.2011.11828420].
- Droeschel, D.; Behnke, S. Efficient Continuous-Time SLAM for 3D Lidar-Based Online Mapping. 2018
 IEEE International Conference on Robotics and Automation (ICRA) 2018, pp. 1–9.
- Vlaminck, M.; Luong, H.; Goeman, W.; Philips, W. 3D Scene Reconstruction Using Omnidirectional Vision
 and LiDAR: A Hybrid Approach. *Sensors* 2016, *16*.
- 437 25. Vlaminck, M.; Luong, H.; Philips, W. Liborg: a lidar-based robot for efficient 3D mapping. APPLICATIONS
 438 OF DIGITAL IMAGE PROCESSING XL. SPIE, 2017, Vol. 10396.
- 439 26. Mellado, N.; Aiger, D.; Mitra, N.J. SUPER 4PCS: Fast Global Pointcloud Registration via Smart Indexing.
 440 *Computer Graphics Forum* 2014.
- Goossens, B.; De Vylder, J.; Philips, W. Quasar: a new heterogeneous programming framework for image
 and video processing algorithms on CPU and GPU. IEEE International Conference on Image Processing
 ICIP. IEEE, 2014, pp. 2183–2185.
- Aiger, D.; Mitra, N.J.; Cohen-Or, D. 4pointss Congruent Sets for Robust Pairwise Surface Registration.
 ACM SIGGRAPH 2008 Papers; ACM: New York, NY, USA, 2008; SIGGRAPH '08, pp. 85:1–85:10.
- 29. Shoemake, K. Animating Rotation with Quaternion Curves. SIGGRAPH Comput. Graph. 1985, 19, 245–254.
- Geiger, A. Are We Ready for Autonomous Driving? The KITTI Vision Benchmark Suite. Proceedings of
 the 2012 IEEE Conference on Computer Vision and Pattern Recognition (CVPR); IEEE Computer Society:
 Washington, DC, USA, 2012; CVPR '12, pp. 3354–3361.
- Wohlkinger, W.; Vincze, M. Ensemble of shape functions for 3D object classification. ROBIO. IEEE, 2011,
 pp. 2987–2992.
- Johnson, A. Spin-Images: A Representation for 3-D Surface Matching. PhD thesis, Carnegie Mellon
 University, Pittsburgh, PA, 1997.
- (c) 2018 by the authors. Submitted to *Sensors* for possible open access publication under the terms and conditions
- of the Creative Commons Attribution (CC BY) license (http://creativecommons.org/licenses/by/4.0/).