



# Point-to-hyperplane RGB-D Pose Estimation: Fusing Photometric and Geometric Measurements

Fernando Israel Ireta Muñoz, Andrew I. Comport

## ► To cite this version:

Fernando Israel Ireta Muñoz, Andrew I. Comport. Point-to-hyperplane RGB-D Pose Estimation: Fusing Photometric and Geometric Measurements. IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS 2016), Oct 2016, Daejeon, South Korea. hal-01324294v2

**HAL Id: hal-01324294**

**<https://hal.science/hal-01324294v2>**

Submitted on 29 Aug 2016

**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.

# Point-to-hyperplane RGB-D Pose Estimation: Fusing Photometric and Geometric Measurements

Fernando I. Ireta Muñoz<sup>1</sup> and Andrew I. Comport<sup>2</sup>

**Abstract**—The objective of this paper is to investigate the problem of how to best combine and fuse color and depth measurements for incremental pose estimation or 3D tracking. Subsequently a framework will be proposed that allows to formulate the problem with a unique measurement vector and not to combine them in an ad-hoc manner. In particular, the full color and depth measurement will be defined as a 4-vector (by combining 3D Euclidean points + image intensities) and an optimal error for pose estimation will be derived from this. As will be shown, this will lead to designing an iterative closest point approach in 4 dimensional space. A kd-tree is used to find the closest point in 4D-space, therefore simultaneously accounting for color and depth. Based on this unified framework a novel Point-to-hyperplane approach will be introduced which has the advantages of classic Point-to-plane ICP but in 4D-space. By doing this it will be shown that there is no longer any need to provide or estimate a scale factor between different measurement types. Consequently, this allows to increase the convergence domain and speed up the alignment, whilst maintaining the robust and accurate properties. Results on both simulated and real environments will be provided along with benchmark comparisons.

## I. INTRODUCTION

Color and depth images acquired from RGB-D sensors are increasingly useful, especially in robotics for computing visual odometry, performing autonomous navigation and reconstructing 3D environments. One of the most fundamental problems is estimating the pose that relates measurements obtained from a moving sensor at different times. Some recent approaches have combined both measurements together in a limited hybrid manner.

The problem of pose estimation from color or depth images have each been individually studied in the computer vision and robotics literature. Classically, color and depth measurements have been used separately in image-based and geometric-based pose estimation. In the case of depth images, the well known Iterative Closest Point (ICP) algorithm prevails [3] and in particular the point-to-plane ICP algorithm is especially efficient and robust [4]. On the other hand, color images have been used to estimate the pose of the camera using direct and dense error functions based on view synthesis [5]. The latter will be referred to here as image-based approaches. It can be noted that feature-based image approaches first extract geometric information from the image before performing estimation on a geometric error. Feature-based approaches are a sub-part of the direct approach and won't be detailed here.

<sup>\*</sup>This work is supported by the European H2020 project: COMANOID, Université Côte d'Azur, CNRS, I3S, France and CONACYT, México.

<sup>1</sup> ireta@i3s.unice.fr <sup>2</sup> Andrew.Comport@cnrs.fr

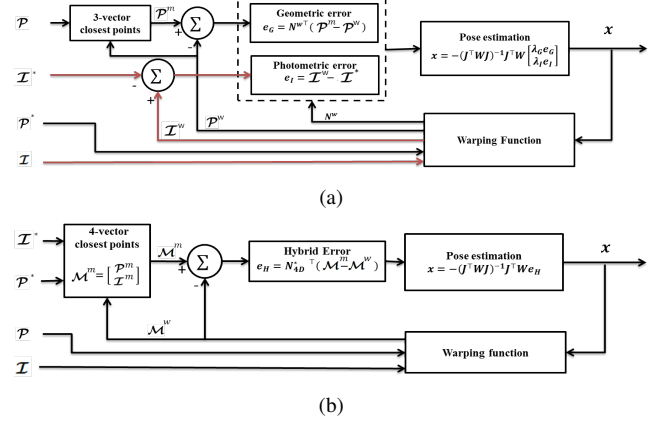


Fig. 1. Two hybrid-based approaches to estimate the unknown pose  $x$  between two sequent RGB-D frames. (a) A direct approach for the color images  $\mathcal{I}$  and  $\mathcal{I}^*$ , and a point-to-plane algorithm for the geometric cloud of points  $\mathcal{P}$  and  $\mathcal{P}^*$  is used. (b) Proposed method: a matching stage considers the 4-vector for minimizing the integrated error with the point-to-hyperplane algorithm, which computes the normals  $N$  in 4D space. The scale factor  $\lambda$  is no longer needed.

Whilst color and depth based pose estimation have been studied separately, similar solutions have been used for both using a non-linear iteratively re-weighted least squares (IRLS) method. The general IRLS pipeline for pose estimation, follows the common strategy across different measurement types:

- 1) Transform/warp the current measurement onto a reference frame using the last pose estimate.
- 2) Find the closest points between the two datasets.
- 3) Determine an error between the two datasets (and robust weights).
- 4) Estimate an incremental update on the pose.
- 5) Repeat to 1. until convergence.

Recently, several strategies have combined color and depth measurements together in different ways and attempt to retain the respective benefits of each. The advantages of using both include increased efficiency, accuracy and robustness. In [22] a recent survey of the real-time performance of these approaches is provided. Image-based approaches alone are dependant on texture in the images to constrain all degrees of freedom. For example, a wall with only horizontal lines would be degenerate. ICP approaches require sufficient geometry and are, for example, degenerate in the case of a movement parallel to a flat wall. In this paper, those approaches which combine depth and color for robust and accurate pose estimation will be referred to as *hybrid* approaches (See Fig. 1).

Amongst the various hybrid methods, those of most interest are those that minimize a photometric and geometric error simultaneously in real-time [12], [13], [20], [21]. The two main differences in the proposed approaches are categorized as:

- How the closest points are determined between different RGB-D measurements.
- How the joint optimization is performed.

The aforementioned hybrid approaches are somewhat ad-hoc because they do not necessarily consider the color and depth simultaneously when computing closest points. Furthermore, in the optimization stage they simply combine the classic ICP and image-based approaches by minimizing both error types simultaneously. This, however, requires the definition or estimation of a tuning parameter  $\lambda$  which weights the respective contribution of each different measurement type.

First, consider a fused version of the closest point search in Step 2 which is required for both ICP and image-based approaches. In the case of ICP alone, the closest points are often obtained by performing a *kd*-tree (k-dimensional) for nearest neighbours search. Alternately, in the image-based approaches the image warping function finds the closest color values by view interpolation (nearest neighbour, bi-linear, bi-cubic,...) directly in image space. Of the recent hybrid approaches [12], [13], [20], [21], each performs the closest point searching separately for both color and depth and no fused information is considered. Methods that consider both color and depth in the closest point matching stage include [11], [15], [14]. The former and later approach use 3 channels of color and differ in the color spaces used while [15] considers only greyscale information. Finding the closest points using both color and depth increases the accuracy of finding the true nearest neighbour, however, this requires an efficient search in 4-space (3D points + intensity) or higher dimensional space.

Now consider the joint optimization problem of the IRLS algorithm that minimizes both a fused ICP and image-based error. The large majority of classic approaches involve simply stacking the two error functions and minimizing the resulting joint error simultaneously [11], [15], [20], [13], [12], [21], [14]. All except [11] perform ICP point-to-plane combined with the image-based approach. The drawback of these approaches is that they require the definition of a tuning parameter  $\lambda$  which weights the respective contribution of each measurement. These methods, then vary in how this tuning parameter is determined. In [11],  $\lambda$  is computed by estimating the ratio between the minimum and maximum values of both, color space and geometric errors and the best value is chosen experimentally in this range. [15] proposes interestingly an adaptive  $\lambda$  which is varied using a sigmoidal function which favors the ICP approach far from the solution and the image-based approach close to the minimum. This has the benefit of faster convergence and more accuracy of the solution. In [20], [13], the scale factor is automatically estimated as the ratio between the Median Absolute Deviations (MAD) of the color error and the

median of the depth error (i.e. their relative robust variance). In [12], [18],  $\lambda$  is estimated by computing the covariance of the residuals for each point individually assuming a *t*-distribution of the error. This improved the convergence rate, however, is computationally expensive to iteratively compute a  $\lambda$  for each pixel.

The aim of this paper is to propose a unified framework for fusing both image-based and ICP strategies for pose estimation at each stage of the IRLS process. As will be shown, this leads to a novel Point-to-hyperplane ICP approach in 4 dimensions (3D + Intensity) which could easily be extended to greater dimensions (for example color RGB). This formulation also naturally leads to a fused closest point search strategy that exploits both color and geometric information simultaneously. The approach used in the paper to find the closest points in 4D space uses a *kd*-tree, however, alternative search strategies could also be used. In practice, and for computational efficiency, the ANN (Approximate Nearest Neighbour) [16] algorithm is used. Furthermore, the *kd*-tree can be built only once from the reference image before the iterative loop, therefore maintaining efficiency.

The paper is organized as follows. Section II briefly explains the classic hybrid approach that jointly minimizes intensities and Point-to-plane ICP. In Section III a novel Point-to-hyperplane approach is introduced. Section IV provides the implementation details common to all the methods that were evaluated. Finally, simulated and real experimental results with benchmarks are presented in Section V.

## II. JOINT METHOD FOR COMBINING GEOMETRIC AND PHOTOMETRIC APPROACHES

Pose estimation from hybrid methods is achieved by fusing the geometric and photometric optimization functions and minimizing the errors simultaneously. The main feature of hybrid methods for estimating the camera poses, is that they constrain the pose estimation better and can converge faster than using the techniques alone.

The pose will be defined here as the homogeneous pose matrix  $\mathbf{T}(\mathbf{x}) \in \mathbb{R}^{4 \times 4}$  which depends on a minimal parameterization of 6 parameters which are defined here as the linear and angular velocity  $\mathbf{x} = [\mathbf{v}, \boldsymbol{\omega}]^\top \in \mathbb{R}^6$ , respectively. The homogeneous transformation matrix can be decomposed into rotational and translational components  $\mathbf{T}(\mathbf{x}) = (\mathbf{R}(\mathbf{x}), \mathbf{t}(\mathbf{x})) \in \text{SE}(3)$ . The relationship between both is given by the exponential map as  $\mathbf{T}(\mathbf{x}) = e^{[\mathbf{x}]_\wedge}$ , with the operator  $[\cdot]_\wedge$  as:

$$[\mathbf{x}]_\wedge = \begin{bmatrix} [\boldsymbol{\omega}]_\times & \mathbf{v} \\ 0 & 0 \end{bmatrix} \quad (1)$$

where  $[\cdot]_\times$  is the skew symmetric matrix operator.

The hybrid approach used to estimate the pose is depicted in Fig. 1(a). It defines an error function that minimizes the joint error between subsequent RGB-D image frames (see [13] for more detail) such as:

$$\mathbf{e}_{H_i} = \rho_i \begin{pmatrix} \lambda \left( \widehat{\mathbf{R}}\mathbf{R}(\mathbf{x})\mathbf{N}_i^* \right)^\top \left( \mathbf{P}_i^m - \Pi_3 \widehat{\mathbf{T}}\mathbf{T}(\mathbf{x})\overline{\mathbf{P}}_i^* \right) \\ \mathbf{I}_i \left( w(\widehat{\mathbf{T}}\mathbf{T}(\mathbf{x}), \mathbf{P}_i^*) \right) - \mathbf{I}_i^* (\mathbf{p}^*) \end{pmatrix} \in \mathbb{R}^4 \quad (2)$$

where the first row of (2) is the Point-to-plane ICP error with the projective data association and the second row is the photometric term. The superscript  $*$  identifies the reference measurements,  $\Pi_3 = [\mathbf{1}, \mathbf{0}] \in \mathbb{R}^{3 \times 4}$  is the projection matrix,  $\mathbf{N}_i^* \in \mathbb{R}^3$  is the surface normal for each homogeneous 3D point  $\mathbf{P}_i^* \in \mathbb{R}^4$ . The closest point  $\mathbf{P}_i^m$  can be obtained by linearly interpolating the warped pixel coordinates into the current depth map as in [12], [13], [20]. The geometric warping function  $w(\cdot)$  projects a reference 3D point  $\mathbf{P}_i^* \in \mathbb{R}^3$  onto the current image plane. The closest image intensity is then found by interpolation of the current intensity function at the warped pixel coordinates to obtain the corresponding intensity as:  $\mathbf{I}_i^w(\mathbf{p}_i^*) = \mathbf{I}_i(\mathbf{p}_i^w) \in \mathbb{Z}$ . The 3D point is computed by the back projection function as  $\mathbf{P}_i = \mathbf{K}^{-1} \bar{\mathbf{p}}_i$   $Z_i = [X_i \ Y_i \ Z_i]^\top \in \mathbb{R}^3$ , where  $\mathbf{K} \in \mathbb{R}^{3 \times 3}$  is the calibration matrix which contains the intrinsic parameters of the camera, and  $Z_i \in \mathbb{R}^+$  is the metric measurement for each pixel coordinate  $\bar{\mathbf{p}}_i = [u_i \ v_i \ 1]^\top \in \mathbb{R}^3$  of the depth image.

The given non-linear error in (2) is minimized iteratively using a Gauss-Newton approach to compute the unknown parameter  $\mathbf{x}$  with increments given by:

$$\mathbf{x} = -(\mathbf{J}^\top \mathbf{W} \mathbf{J})^{-1} \mathbf{J}^\top \mathbf{W} \begin{bmatrix} \lambda \mathbf{e}_G \\ \mathbf{e}_I \end{bmatrix} \quad (3)$$

where  $\mathbf{J} = [\mathbf{J}_I \ \mathbf{J}_G]^\top$  represents the stacked Jacobian matrices obtained by derivation of the stacked photometric and geometric error functions ( $\mathbf{e}_I$  and  $\mathbf{e}_G$  respectively), and the weight matrix  $\mathbf{W}$  contains the weights  $\rho_i$  associated to each set of coordinates obtained by M-estimation [9]. The photometric Jacobian  $\mathbf{J}_I$  is computed using the efficient second order minimization method (ESM) [2]. The pose estimate  $\mathbf{T}(\mathbf{x})$  is computed at each iteration and is updated incrementally as  $\hat{\mathbf{T}} \leftarrow \hat{\mathbf{T}} \mathbf{T}(\mathbf{x})$  until convergence.

The parameter  $\lambda$  is a constant that scales the relative error distributions. As mentioned in the introduction, many methods have been proposed to estimate this parameter ranging from manual tuning to more complex estimation. Manually fixing  $\lambda$  is not optimal nor efficient and estimating the parameter requires extra computational cost. For the purposes of this paper,  $\lambda$  has been estimated as in [20]. In the following section, is going to be seen that  $\lambda$  is not required if we consider a Point-to-hyperplane approach.

### III. POINT-TO-HYPERPLANE METHOD

As mentioned in the introduction, the objective of this paper is to perform both closest point matching and minimization using a 4-vector containing color and depth. Since 4D space has an additional degree of freedom, the normal obtained for the 3D point-to-plane method will be orthogonal to a surface in 4D which spans both geometry and color. This surface will be referred in this paper as *hyperplane*. The 4-vector is defined as  $\mathbf{M}_i = [\mathbf{P}_i^\top \ \mathbf{I}_i]^\top \in \mathbb{R}^4$ , where the 3D Euclidean point  $\mathbf{P}_i$  is fused with its associated greyscale intensity  $\mathbf{I}_i$  in a single measurement vector.

Two measurement vectors obtained at different views of the same scene are generally not in correspondence. The fused error can then be defined as a 4D IRLS problem

between two point clouds. The hybrid error function (case 1) is defined such as:

$$\begin{aligned} \mathbf{e}_{H_i} &= \rho_i (\mathbf{N}_i^{*\top} (\mathbf{M}_i^* - \mathbf{M}_i^m)) \in \mathbb{R}^4 \\ &= \rho_i \left( \mathbf{N}_i^{*\top} \begin{pmatrix} \mathbf{P}_i^* - \Pi_3 \hat{\mathbf{T}} \mathbf{T}(\mathbf{x}) \bar{\mathbf{P}}_i^m \\ \mathbf{I}_i^* - \mathbf{I}_i(w(\hat{\mathbf{T}} \mathbf{T}(\mathbf{x}), \mathbf{P}_i^*)) \end{pmatrix} \right) \in \mathbb{R}^4 \end{aligned} \quad (4)$$

where  $\mathbf{M}_i^*$  is the reference 4D point,  $\mathbf{M}_i^m$  corresponds to the warped closest points to image according to the unknown transformation  $\hat{\mathbf{T}} \mathbf{T}(\mathbf{x})$ . This is similar to (2) except that the normal is computed in 4D and the closest points can be determined in 4D. Also, note that the current measurements in (4) are all warped to the reference frame so that it is no longer necessary to rotate the normals as in (2). This formulation seems better than in [13] and the computational time is negligible.

Alternately, it is also possible to find the closest points in the current 4-vector (case 2) that solve the optimization problem as:

$$\mathbf{e}_{H_i} = \rho_i (\mathbf{N}_i^{m\top} (\mathbf{M}_i^m - \mathbf{M}_i^w)) \in \mathbb{R}^4 \quad (5)$$

where  $\mathbf{M}_w$  is the warped 4-vector. Solutions in both cases will be considered, but first lets consider the 4D normal.

Mathematically, at least four 4D-points are needed to compute three vectors that will be used to perform a three-way cross product to obtain the 4D-normal. The cross product does not, however, account for the uncertainty of the points in its computation. Instead, a Principal Component Analysis (PCA) can be computed on the covariance matrix of the nearest neighbours surrounding each point. In that case the smallest eigenvalue corresponds to the surface normal [17]. An alternative solution to compute normals fast and accurately is given in [1]. Considering that  $\mathbf{N}_i^m$  is directly obtained from  $\mathbf{N}_i^*$ , it is only necessary to compute the normals once for the reference image as in the case of the inverse compositional algorithm for image-based registration. The covariance matrix, used in computing the 4D-normal, allows to project and weight the errors such that the Point-to-hyperplane error is invariant to any tuning parameter  $\lambda$  (See [10] for the proof).

Therefore, the error functions (4) or (5) can be minimized iteratively without scaling the geometric and photometric error such as:

$$\mathbf{x} = -(\mathbf{J}^\top \mathbf{W} \mathbf{J})^{-1} \mathbf{J}^\top \mathbf{W} \mathbf{e}_H \quad (6)$$

where  $\mathbf{e}_H$  is the stacked hybrid error  $\mathbf{N}_i^{*\top} (\mathbf{M}_i^* - \mathbf{M}_i^m)$  or  $\mathbf{N}_i^{m\top} (\mathbf{M}_i^m - \mathbf{M}_i^w)$ , and the weights  $\rho_i$  are stacked in  $\mathbf{W}$ .

As presented above in (4) and (5) it is possible to find the closest points either in the reference frame or the current frame. Consider first case 1. Computing in the reference frame has the advantage that several parameters can be pre-computed only once and not at each iteration. In that case the reference normals can be precomputed along with a quick search strategy for finding closest points such as a kd-tree. Consider now case 2. When finding the closest points in the current image it is not possible to recompute the kd-tree for each new image because it is computationally too expensive.

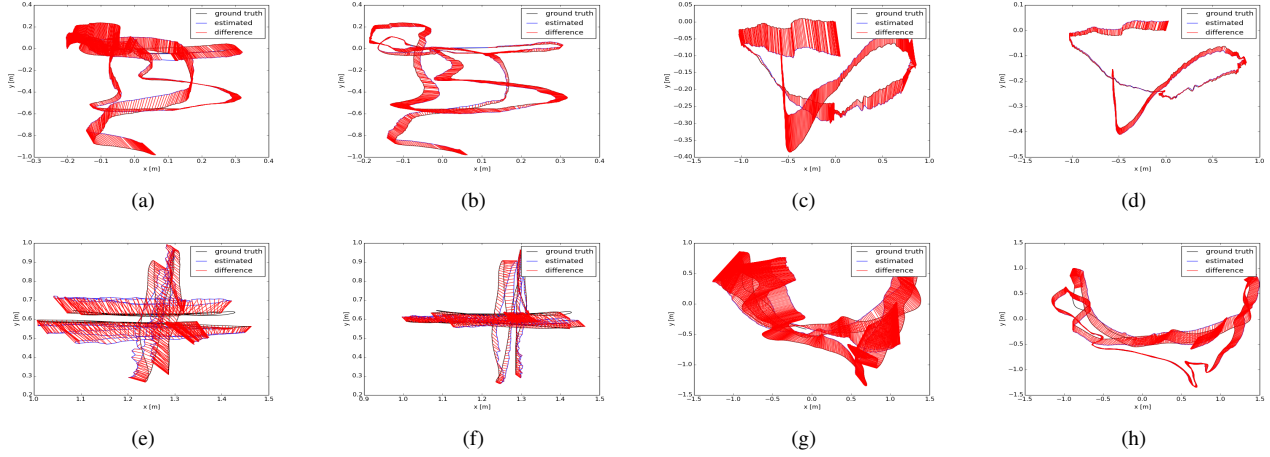


Fig. 2. Examples of the Absolute Trajectory Error evaluation, the first and third column show the results obtained by methods that combine the direct approach with the geometric Point-to-plane method. The second and fourth column for the point-to-hyperplane method. The trajectories presented here are obtained for the simulated (a)(b) lvr/traj0, (c)(d) lvr/traj2, [7] and the real (e)(f) fr1/xyz, (g)(h) fr1/room [19] sequences. More sequences are shown in [10].

In this case it is possible to consider approximating the closest point by simply search for the closest point in the image (as is done in (2)). In this case nearest neighbour, bilinear or bi-cubic interpolation can be performed [12], [13], [20], [21]. In the case where no depth is provided by the RGB-D sensor for a pixel, it is rejected, unless it is possible to interpolate a depth value from neighboring pixels.

Another strategy often used in the ICP literature (see for example [6], [8]), is to compute closest point matching only in the first iteration of the IRLS minimization loop. This allows to avoid too much computational complexity while obtaining the benefits of finding the closest points. In this paper a classic  $kd$ -tree with an ANN algorithm is considered so that large displacements will converge to a solution. To estimate the closest points, the optimized search function of the FLANN library [16] is employed to find the true nearest neighbours in  $4D^1$ . In a local tracking situation, faster projective data association techniques can be employed.

#### IV. IMPLEMENTATION DETAILS

In this section, some parameters considered for the experiments will be established. The experiments were done for real and synthetic RGB-D greyscale images in MATLAB. A multi-resolution pyramid was used to improve the computational efficiency. The images were warped at the second level of the pyramid (resolution  $160 \times 120$ ).

There are two convergence criteria that stop the iterative loop for real and simulated experiments. The first one is a maximum number of iterations, which is established as 200, and the norm of the estimated rotation and translation. If the transformation matrix  $T(x)$  gets closer to the identity matrix, then the iterative loop stops. The parameters used to determine this second break are  $norm(\mathbf{x}_R) < 1 \times 10^{-6}$  for rotation and  $norm(\mathbf{x}_t) < 1 \times 10^{-5}$  for translation.

<sup>1</sup>The legend NN4D is employed throughout the paper to indicate that the closest points were estimated in the first iteration only by finding the nearest neighbours with a  $kd$ -tree in the 4-vector.

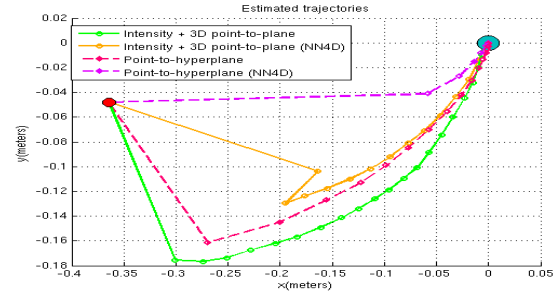


Fig. 4. Example of 4 estimated camera trajectories between a pair of images with a random pose in a simulated environment. The green and red dot indicate the initial and the final pose. The Point-to-hyperplane methods improve the Point-to-plane + direct methods, obtaining more direct trajectories when NN4D are estimated. 1000 synthesized images with a random pose were equally tested, obtaining a similar performance.

To reject outliers, M-estimators were employed. They are more general because they permit the use of different minimization functions not necessarily corresponding to normally distributed data. In this paper, the Huber influence function was used for this purpose.

The normal in 4D is adjusted to a  $3 \times 3$  window that perform the PCA algorithm to find the covariance matrix of the fused error, leading to find the smallest eigenvalue which corresponds to the normal parameter.

All the experiments were validated on a workstation with Ubuntu 14.04, Intel Core i7-4770K and 16 GB RAM.

#### V. RESULTS

During the experimental part, it was seen that the parameter  $\lambda$  does not change the performance or accuracy of the pose estimation when the Point-to-hyperplane approach is used, but it improves the performance of other hybrid approaches if it is well estimated. An introduction to the improvement can be seen in Fig. 4, where the tracking trajectories estimated by 4 strategies are shown.

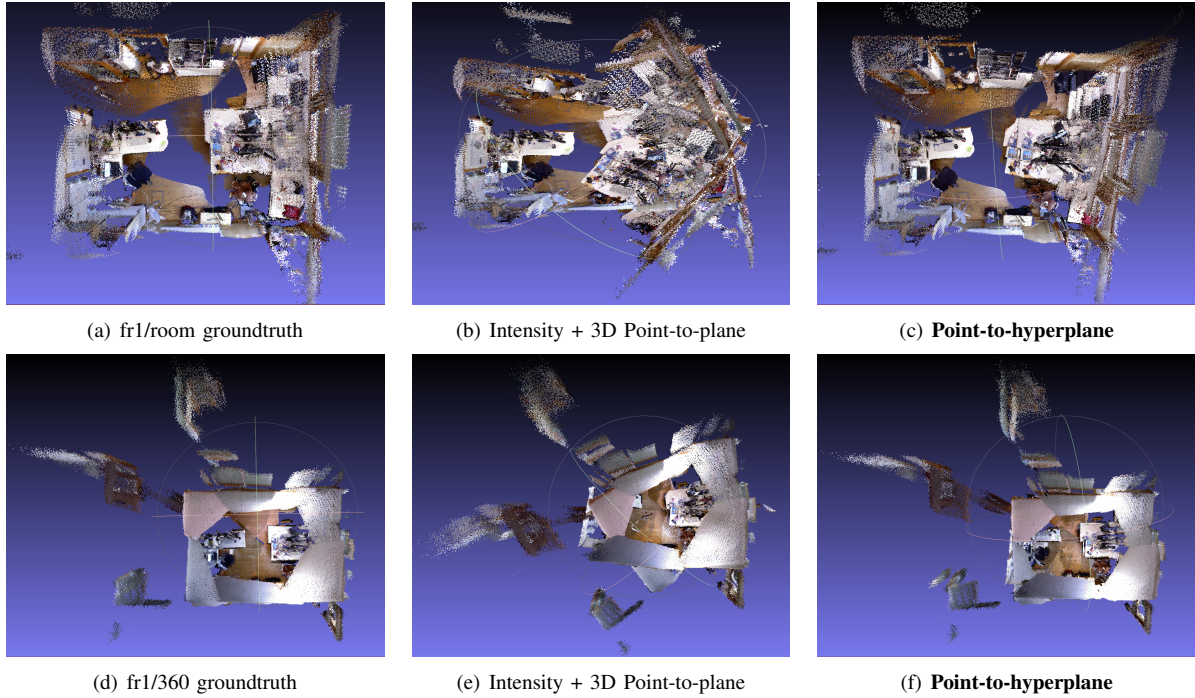


Fig. 3. 3D reconstruction of sequences fr1/room and fr1/360 (first and second row, respectively). In the first column the groundtruth obtained by an external motion capture system is shown, the column in the middle shown the results of the direct approach + 3D point to plane algorithm and the last column shown the result of the Point-to-hyperplane method. This difficult 360 degree sequence with motion blur clearly shows that the proposed method can achieve more robust estimations. The  $\lambda$  parameter was set initially using the automatic approach from [20] and was not tuned individually for each frame of the sequence. Therefore, the method could be improved with strategies as loop closure detection algorithms and keyframe detectors.

TABLE I

RELATIVE POSE ERROR (RPE) AND ABSOLUTE TRAJECTORY ERROR (ATE) FOR THE SIMULATED AND REAL DATASET [7], [19]. IT CAN BE SEEN THAT THE POINT-TO-HYPERPLANE METHODS (3 & 4) IMPROVE THE HYBRID METHODS THAT COMBINE THE DIRECT APPROACH AND THE GEOMETRIC POINT-TO-PLANE (1 & 2) IN THE MAJORITY OF DATASET FOR THE RPE TRANSLATIONAL EVALUATION AND IN ALL DATASET FOR RPE ROTATIONAL EVALUATION. NOTE THAT THE ONLY DIFFERENCE BETWEEN 1 & 2 AND 3 & 4 IS THE KD-TREE MATCH ON THE FIRST ITERATION, WHICH ONLY AFFECTS COMPUTATION TIME AND NOT ACCURACY. THE REMAINING ITERATIONS USE THE SAME ERROR FUNCTION.

Sequence	Method	RPE translational (m)			RPE rotational (deg)			ATE (m)		
		RMSE	MEAN	STD	RMSE	MEAN	STD	RMSE	MEAN	STD
fr1/xyz	1 & 2	0.033	0.030	0.014	2.025	1.741	1.034	0.095	0.087	0.039
	3 & 4	<b>0.021</b>	<b>0.019</b>	<b>0.008</b>	<b>1.106</b>	<b>0.998</b>	<b>0.477</b>	<b>0.045</b>	<b>0.038</b>	<b>0.024</b>
fr1/rpy	1 & 2	0.062	0.050	0.037	3.161	2.887	1.288	0.136	0.115	0.072
	3 & 4	<b>0.038</b>	<b>0.032</b>	<b>0.020</b>	<b>2.820</b>	<b>2.652</b>	<b>0.959</b>	<b>0.035</b>	<b>0.032</b>	<b>0.015</b>
fr1/360	1 & 2	<b>0.146</b>	0.118	<b>0.086</b>	4.171	3.844	1.621	0.520	0.484	0.188
	3 & 4	0.152	<b>0.114</b>	0.100	<b>3.159</b>	<b>2.859</b>	<b>1.343</b>	<b>0.322</b>	<b>0.296</b>	<b>0.125</b>
fr1/room	1 & 2	0.076	0.060	0.048	3.285	2.912	1.520	0.434	0.404	0.158
	3 & 4	<b>0.056</b>	<b>0.047</b>	<b>0.030</b>	<b>2.673</b>	<b>2.329</b>	<b>1.313</b>	<b>0.174</b>	<b>0.152</b>	<b>0.086</b>
fr1/desk	1 & 2	0.047	0.039	0.027	2.826	2.503	1.312	0.108	0.104	0.029
	3 & 4	<b>0.044</b>	<b>0.036</b>	<b>0.025</b>	<b>2.309</b>	<b>2.027</b>	<b>1.106</b>	<b>0.071</b>	<b>0.067</b>	<b>0.023</b>
fr1/desk2	1 & 2	<b>0.058</b>	<b>0.051</b>	<b>0.027</b>	3.483	3.026	1.725	0.189	0.174	0.075
	3 & 4	0.060	<b>0.051</b>	0.031	<b>3.026</b>	<b>2.641</b>	<b>1.478</b>	<b>0.133</b>	<b>0.116</b>	<b>0.065</b>
fr1/floor	1 & 2	0.094	<b>0.038</b>	0.086	4.660	1.953	4.231	0.772	0.666	0.391
	3 & 4	<b>0.080</b>	0.051	<b>0.062</b>	<b>3.909</b>	<b>1.915</b>	<b>3.408</b>	<b>0.473</b>	<b>0.405</b>	<b>0.244</b>
fr1/plant	1 & 2	0.106	0.067	0.082	3.941	3.223	2.268	0.324	0.296	0.132
	3 & 4	<b>0.055</b>	<b>0.043</b>	<b>0.034</b>	<b>2.130</b>	<b>1.947</b>	<b>0.864</b>	<b>0.101</b>	<b>0.093</b>	<b>0.037</b>
fr1/teddy	1 & 2	0.096	0.081	0.051	3.410	3.021	1.583	0.615	0.553	0.271
	3 & 4	<b>0.070</b>	<b>0.056</b>	<b>0.043</b>	<b>2.287</b>	<b>1.954</b>	<b>1.187</b>	<b>0.169</b>	<b>0.158</b>	<b>0.059</b>
lvr/traj0	1 & 2	<b>0.001</b>	<b>0.001</b>	<b>0.001</b>	0.044	0.035	<b>0.027</b>	0.128	0.114	0.057
	3 & 4	0.002	<b>0.001</b>	0.002	<b>0.042</b>	<b>0.026</b>	0.033	<b>0.050</b>	<b>0.046</b>	<b>0.019</b>
lvr/traj1	1 & 2	0.002	<b>0.001</b>	<b>0.001</b>	0.048	0.041	0.024	0.114	0.104	0.046
	3 & 4	<b>0.001</b>	<b>0.001</b>	<b>0.001</b>	<b>0.021</b>	<b>0.017</b>	<b>0.013</b>	<b>0.041</b>	<b>0.032</b>	<b>0.026</b>
lvr/traj2	1 & 2	0.002	<b>0.001</b>	<b>0.001</b>	0.044	0.039	0.021	0.074	0.067	0.030
	3 & 4	<b>0.001</b>	<b>0.001</b>	<b>0.001</b>	<b>0.024</b>	<b>0.019</b>	<b>0.014</b>	<b>0.039</b>	<b>0.036</b>	<b>0.016</b>
lvr/traj3	1 & 2	0.002	<b>0.001</b>	<b>0.001</b>	0.070	0.053	0.045	0.218	0.202	0.082
	3 & 4	<b>0.001</b>	<b>0.001</b>	<b>0.001</b>	<b>0.044</b>	<b>0.027</b>	<b>0.035</b>	<b>0.080</b>	<b>0.066</b>	<b>0.045</b>



TABLE II

AVERAGES IN TIME AND IN NUMBER OF ITERATIONS UNTIL CONVERGENCE FOR 1000 SYNTHESIZED IMAGES AT RANDOM POSES.

Method	# Iterations	Time (sec)
1) Intensity + point-to-plane	65.65	0.61
2) Intensity + point-to-plane (NN4D)	<u>34.23</u>	<u>0.35</u>
3) Point-to-hyperplane	53.24	0.56
4) <b>Point-to-hyperplane (NN4D)</b>	<b>12.83</b>	<b>0.16</b>

1) *Simulated environment - Tracking*: A random transformation is applied to a reference RGB-D image. The new synthesized image is considered as the current image and the methods finds the alignment between each new image and the reference. The motivation for using synthetic data is that the generated images provide a groundtruth for evaluation, since the correspondences between the transformed views are known. The averages shown in Table II demonstrates that the Point-to-hyperplane method improves the number of iterations and computational time<sup>2</sup>.

2) *Real environment - Visual odometry*: The methods were evaluated on the ICL-NUIM RGB-D benchmark dataset [7] and on full benchmarks *freiburg 1* sequences from TUM [19]. Therefore, in order to evaluate the estimated trajectories frame-to-frame, the online tool was used with the default settings to evaluate the ATE (Absolute Trajectory Error) and RPE (Relative Pose Error), which are compared alongside an accurate groundtruth trajectory in Table I. Some examples of the improvement are shown in Fig. 2. A demonstration of the performance of the proposed method can be seen in the video attachment of this paper.

## VI. CONCLUSION

A novel Point-to-hyperplane strategy was proposed based on a 4D-vector. Two main advantages of this unified framework are underlined. First, it is shown that the sensor pose can be estimated by minimizing the combined error without computing a scale parameter between them. Second, it is shown that nearest neighbour techniques can be used in 4D-space to improve the convergence rate and domain for IRLS pose estimation. Experimental results and analysis are provided which compare two variants of the new Point-to-hyperplane approach with the classic hybrid approach. The results show improved computation time and faster convergence on well known benchmarks. Future work will be dedicated to testing the new approach in a simultaneous localization and mapping context which will provide for interesting comparisons with the map reconstruction. The framework would also apply directly to volumetric representations by modifying the projective error function accordingly. We will also extend the strategy to higher dimensions by using color and other measurement to increase robustness.

<sup>2</sup>The time shown does not consider the computation of normals nor construction of the kd-tree. A more detailed computational analysis can be found in [10].

## REFERENCES

- [1] H. Badino, D. Huber, Y. Park, and T. Kanade. Fast and accurate computation of surface normals from range images. In *ICRA 2011*, pages 3084–3091, May 2011.
- [2] Selim Benhimane and E. Malis. Real-time image-based tracking of planes using efficient second-order minimization. In *Intelligent Robots and Systems, 2004. (IROS 2004). Proceedings. 2004 IEEE/RSJ International Conference on*, volume 1, pages 943–948 vol.1, Sept 2004.
- [3] P.J. Besl and Neil D. McKay. A method for registration of 3-d shapes. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 14(2):239–256, Feb 1992.
- [4] Yang Chen and Gérard Medioni. Object modelling by registration of multiple range images. *Image Vision Comput.*, 10(3):145–155, April 1992.
- [5] Andrew I Comport, Ezio Malis, and Patrick Rives. Accurate quadric focal tracking for robust 3d visual odometry. In *Robotics and Automation, 2007 IEEE International Conference on*, pages 40–45. IEEE, 2007.
- [6] S. Druon, M.J. Aldon, and A. Crosnier. Color constrained icp for registration of large unstructured 3d color data sets. In *Information Acquisition, 2006 IEEE International Conference on*, pages 249–255, Aug 2006.
- [7] A. Handa, T. Whelan, J. McDonald, and A.J. Davison. A benchmark for rgb-d visual odometry, 3d reconstruction and slam. In *Robotics and Automation (ICRA), 2014 IEEE International Conference on*, pages 1524–1531, May 2014.
- [8] Peter Henry, Michael Krainin, Evan Herbst, Xiaofeng Ren, and Dieter Fox. *Experimental Robotics: The 12th International Symposium on Experimental Robotics*, chapter RGB-D Mapping: Using Depth Cameras for Dense 3D Modeling of Indoor Environments, pages 477–491. Springer Berlin Heidelberg, Berlin, Heidelberg, 2014.
- [9] P.J. Huber, J. Wiley, and W. InterScience. *Robust statistics*. Wiley New York, 1981.
- [10] Fernando I. Ireta Muñoz and Andrew I. Comport. A proof that fusing measurements using point-to-hyperplane registration is invariant to relative scale. In *IEEE International Conference on Multisensor Fusion and Integration for Intelligent Systems*, Baden-Baden, Germany, 2016.
- [11] Andrew Edie Johnson and Bing Kang. "registration and integration of textured 3d data ". *Image and Vision Computing*, 17(2):135 – 147, 1999.
- [12] C. Kerl, J. Sturm, and D. Cremers. Dense visual slam for rgb-d cameras. In *IROS*, 2013.
- [13] M. Meilland and A.I. Comport. On unifying key-frame and voxel-based dense visual SLAM at large scales. In *International Conference on Intelligent Robots and Systems*, Tokyo, Japan, 3-8 November 2013. IEEE/RSJ.
- [14] J. Pauli Michael Korn, M. Holzkoth. Color supported generalized-icp. In *International Conference on Computer Vision Theory and Applications*, 2014.
- [15] L. Morency and T. Darrell. Stereo tracking using icp and normal flow constraint. In *Pattern Recognition, 2002. Proceedings. 16th International Conference on*, volume 4, pages 367–372 vol.4, 2002.
- [16] Marius Muja and David G. Lowe. Scalable nearest neighbor algorithms for high dimensional data. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 36, 2014.
- [17] A. Segal, D. Haehnel, and S. Thrun. Generalized-icp. In *Proceedings of Robotics: Science and Systems*, Seattle, USA, June 2009.
- [18] F. Steinbruecker, C. Kerl, J. Sturm, and D. Cremers. Large-scale multi-resolution surface reconstruction from rgb-d sequences. In *ICCV*, Sydney, Australia, 2013.
- [19] J. Sturm, N. Engelhard, F. Endres, W. Burgard, and D. Cremers. A benchmark for the evaluation of rgb-d slam systems. In *Proc. of the International Conference on Intelligent Robot Systems (IROS)*, Oct. 2012.
- [20] T.M. Tykkälä, C. Audras, and A.I Comport. Direct Iterative Closest Point for Real-time Visual Odometry. In *The Second international Workshop on Computer Vision in Vehicle Technology: From Earth to Mars, ICCV*, Barcelona, Spain, November 6-13 2011.
- [21] T. Whelan, H. Johannsson, M. Kaess, J.J. Leonard, and J. McDonald. Robust real-time visual odometry for dense rgb-d mapping. In *Robotics and Automation (ICRA), 2013*, pages 5724–5731, May 2013.
- [22] Qian-Yi Zhou and Vladlen Koltun. Color map optimization for 3d reconstruction with consumer depth cameras. *ACM Trans. Graph.*, 33(4):155:1–155:10, July 2014.