

Viewpoint-free Video Synthesis with an Integrated 4D System

Csaba Benedek, Zsolt Jankó, Attila Börös, Iván Eichhardt, Dmitry Chetverikov, and Tamás Szirányi
Institute for Computer Science and Control, Hungarian Academy of Sciences (MTA SZTAKI)
{firstname.lastname}@sztaki.mta.hu
<http://web.eee.sztaki.hu/i4d>

Abstract

In this paper, we introduce a complex approach on 4D reconstruction of dynamic scenarios containing multiple walking pedestrians. The input of the process is a point cloud sequence recorded by a rotating multi-beam Lidar sensor; which monitors the scene from a fixed position. The output is a geometrically reconstructed and textured scene containing moving 4D people models, which can follow in real time the trajectories of the walking pedestrians observed on the Lidar data flow. Our implemented system consists of four main steps. First, we separate foreground and background regions in each point cloud frame of the sequence by a robust probabilistic approach. Second, we perform moving pedestrian detection and tracking, so that among the point cloud regions classified as foreground, we separate the different objects, and assign the corresponding people positions to each other over the consecutive frames of the Lidar measurement sequence. Third, we geometrically reconstruct the ground, walls and further objects of the background scene, and texture the obtained models with photos taken from the scene. Fourth we insert into the scene textured 4D models of moving pedestrians which were preliminary created in a special 4D reconstruction studio. Finally, we integrate the system elements in a joint dynamic scene model and visualize the 4D scenario.

Categories and Subject Descriptors (according to ACM CCS): I.4.5 [Computer vision]: Reconstruction

1. Introduction

Recently, an internal project called Integrated 4D (or just i4D) has been launched by two units of the Institute for Computer Science and Control of the Hungarian Academy of Sciences (MTA SZTAKI). The name of the project refers to an unconventional attempt to combine two very different sources of spatio-temporal information, namely, a LIDAR and a 4D reconstruction studio. The main motivation for the integration of the two types of data is our desire to measure and represent the visual world at different levels of detail.

A LIDAR sensor provides a global description of a dynamic outdoor scene in the form of a time-varying 3D point cloud. The latter is used to separate moving objects from static environment and obtain a 3D model of the environment. A 4D studio builds a detailed dynamic model of an actor (typically, a person) moving in the studio. By integrating the two sources of data, one can modify the model of the scene and populate it with the avatars created in the studio. In this paper, we report on the current state of the ongoing i4D project and describe all major processing steps of the integrated system, from the acquisition of the raw data (point

clouds and videos) to the creation and visualisation of an augmented spatio-temporal model of the scene.

LIDAR sensors have been traditionally used in applications such as road extraction in urban areas⁹, vehicle safety and environment recognition¹⁵, airborne data processing and digital terrain modelling^{5,14}, measurement of trees in a forest¹⁷, and modelling of buildings¹⁹ and other constructions. LIDAR data have been integrated with high resolution imagery⁹ and fused with multispectral data¹⁷.

To our best knowledge, there has been no previous attempt to integrate LIDAR data with the output of a 4D reconstruction studio before¹. A typical 4D studio is a green or blue “box” equipped with multiple synchronised, calibrated video cameras. The video streams are used to create a dynamic model of an actor in real time or offline. The degree of realism in shape and appearance varies depending on the approach and the facilities, but the motion of the model obtained in a 4D studio is usually more realistic than that of an artificially created CAD model. Recently, we have built at MTA SZTAKI a 4D studio operating offline¹⁰ and in real time⁸. Sec. 3 gives a brief description of our 4D studio. The

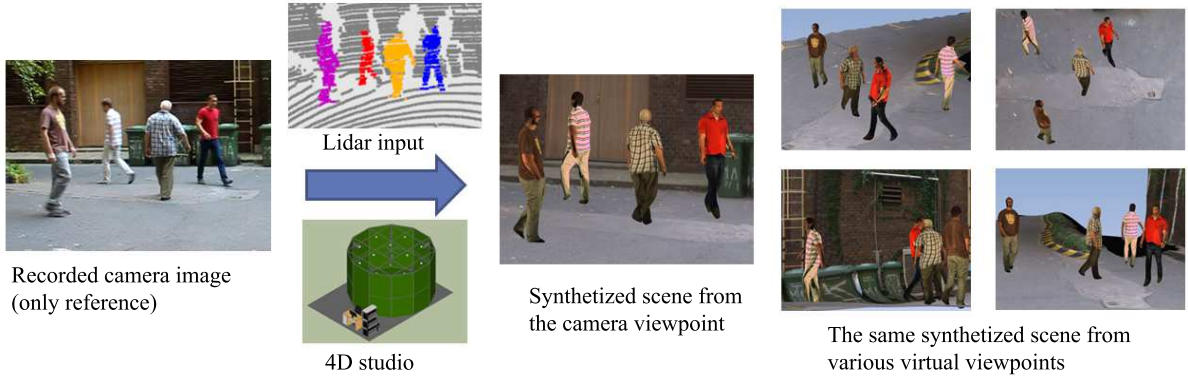


Figure 1: Demonstration of the integrated 4D reconstruction system.

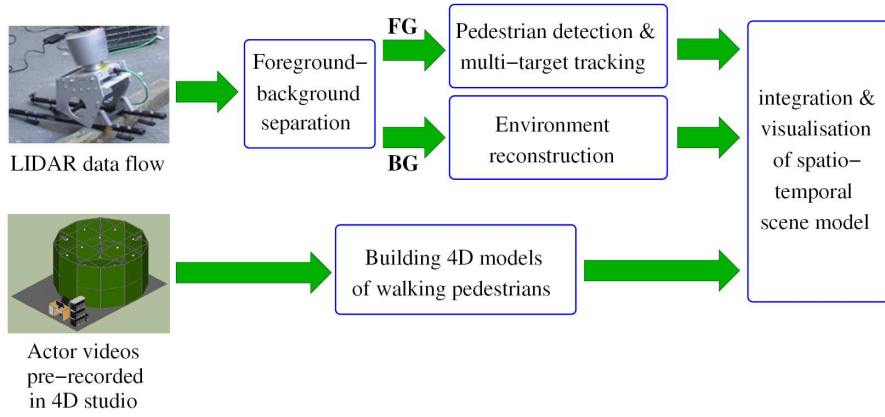


Figure 2: Flowchart of the proposed system. *BG* is background, *FG* foreground.

reader is referred to paper ¹⁰ for a survey of advanced 4D studios in Europe and the USA, and a discussion of their applications in game production, film industry, and other areas.

2. LIDAR data processing

In this section, we present a hybrid method for dense foreground-background point labelling in a point cloud obtained by a Velodyne HDL-64E RMB-LIDAR device that monitors the scene from a fixed position. The method solves the computationally critical spatial filtering tasks applying an MRF model in the 2D range image domain. The ambiguities of the point-to-pixel mapping are handled by joint consideration of the true 3D positions and the 2D labels. Then, we execute detection and tracking of moving pedestrians for the foreground points. Next, we transform the background point cloud into a polygon mesh while maintaining the information about individual objects such as ground, walls, and trees. Finally, the models of the environment objects are

manually textured using photos taken in the scene. Below, we describe these steps in more detail.

2.1. Foreground-background separation

The rotating multi-beam LIDAR device records 360°-view-angle range data sequences of irregular point clouds. Examples of measured point clouds will be shown later in this paper. To separate dynamic foreground from static background in a range data sequence, we apply a probabilistic approach ².

To ensure real-time operation, we project the irregular point cloud to a cylinder surface yielding a depth image on a regular lattice, and perform the segmentation in the 2D range image domain. A part of a range image showing several pedestrians is demonstrated in Fig. 3. Spurious effects are caused by the quantisation error of the discretised view angle, the non-linear position corrections of sensor calibration, and the background flickering, e.g., due to vegetation motion.

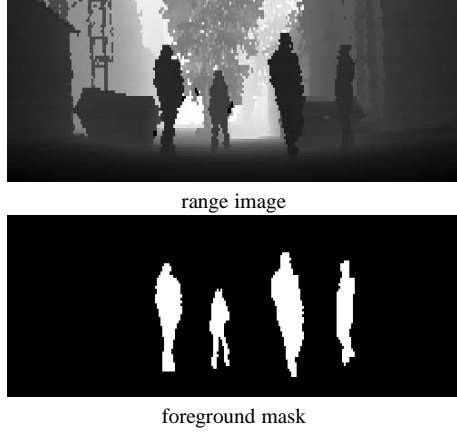


Figure 3: Example of foreground-background segmentation.

One can model the dynamic range image as a Mixture of Gaussians and update the parameters similarly to the standard approach¹⁸. This provides a segmentation of the point cloud which is quite noisy because of the spurious effects. These effects are significantly decreased by the dynamic MRF model² that describes the background and foreground classes by both spatial and temporal features. The model is defined in the range image space. The 2D image segmentation is followed by a 3D point classification step to resolve the ambiguities of the 3D-2D mapping. Using a spatial foreground model, we remove a large part of the irrelevant background motion which is mainly caused by moving tree crowns. Fig. 3 shows an example of foreground segmentation.

2.2. Pedestrian detection and multi-target tracking

In this section, we present the pedestrian tracking module of the system. The input of the module step is a point cloud sequence, where each point is marked with a segmentation label of foreground or background. The output consists of clusters of foreground regions so that the points corresponding to the same person receive the same label over the sequence. We also generate a 2D foot point trajectory of each pedestrian to be used by the 4D scene reconstruction module.

First, the point cloud regions classified as foreground are clustered to obtain separate blobs for each moving person. We fit a regular lattice to the ground plane and project foreground regions onto this lattice. Morphological filters are applied in the image plane to obtain spatially connected blobs for different persons. Then we extract appropriately sized connected components that satisfy area constraints determined by lower and higher thresholds.

This procedure is illustrated in Fig. 4. The centre of each extracted blob is considered as a candidate for foot position

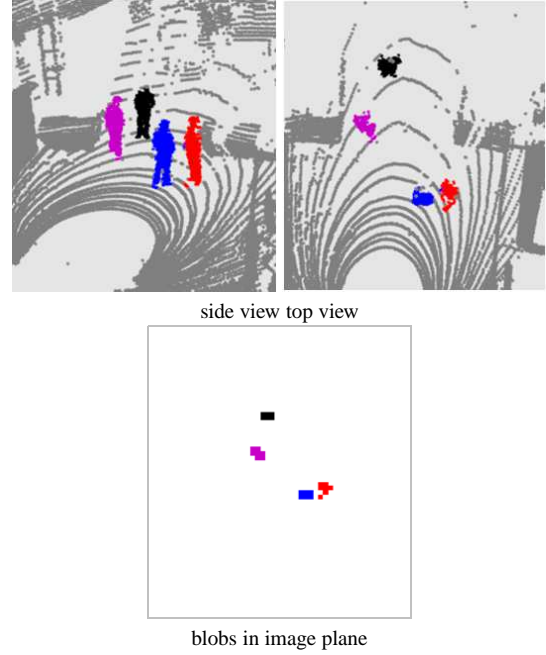


Figure 4: Illustration of pedestrian separation.

in the ground plane. Connected pedestrian shapes may be merged into one blob, while blobs of partially occluded persons may be missed or broken into several parts. Instead of proposing various heuristic rules to eliminate these artefacts at the level of the individual time frames, we developed a robust multi-tracking module which efficiently handles the problems at the sequence level.

Our multi-tracking algorithm receives the measured ground plane positions and for each frame iterates three basic operations, namely, data assignment, Kalman filter correction and Kalman filter prediction. The assignment operation assigns the candidate positions to objects, then the object positions are corrected and, finally, predictions for the subsequent positions are made and fed back to the assignment procedure. The algorithm can handle false positives as well as tracks starting and terminating within a sequence. Temporary track discontinuities are bridged in a post-processing step, while short false tracks are removed based on their length.

The tracker module provides a set of pedestrian trajectories, which are 2D foot centre point sequences in the ground plane. To determine the points corresponding to each pedestrian in a selected frame, the connected foot blobs around a given trajectory point should be vertically back-projected to the 3D point cloud. A result of tracking is demonstrated in Fig. 5 that shows two segmented point cloud frames from a measurement sequence in a courtyard. It also shows the video frames taken in parallel as reference. One can observe

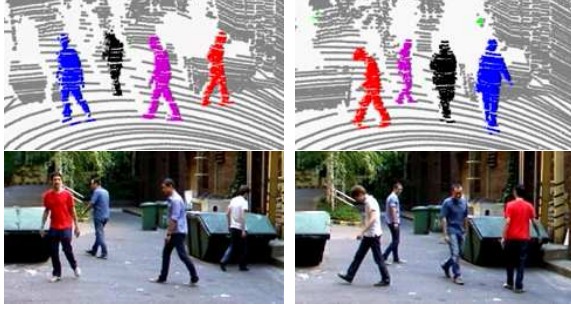


Figure 5: Example of pedestrian tracking in a LIDAR sequence. Top row: point clusters whose colours identify the tracked persons. Bottom row: corresponding video frames displayed for verification.

that during the tracking the point cluster of a pedestrian preserves its colour.

2.3. Environment reconstruction

In this section, we describe our method for static environment reconstruction. First we accumulate the background points of the LIDAR sequence collected over several frames, which results in a dense point cloud that represents the ground, walls, trees, and other background objects. Assuming that the ground is reasonably flat and horizontal, we fit an optimal plane to this point cloud using the robust RANSAC⁷ algorithm that treats all other objects as outliers. Points close to this plane are considered as ground points in the following. For vegetation detection and removal, we have developed an algorithm, which calculates a statistical feature for each point in the merged point cloud based on the distance and irregularity of its neighbors, and also exploits the intensity channel which is an additional indicator of vegetation, which reflects the laser beam with a lower intensity. The remaining points are then projected vertically to the obtained ground plane, where projections of wall points form straight lines that are extracted by the Hough transform⁶. Applying the Ball-Pivoting algorithm³ to the 3D points that project to a straight line, we create a polygon mesh of a wall.

In the reconstruction phase, static background objects of the scene, such as trees, containers or parking cars are replaced with 3D models obtained from Google’s 3D Warehouse. The recognition of these objects from the point cloud is currently done manually, and we are now working on the automation of this step. For example, one can adopt here the machine learning based approach of¹², which extracts various object level descriptors for point cloud blobs representing the detected objects, while to obtain similar representations of the training models from the 3D Warehouse, they perform ray casting on the models to generate point clouds, finally the classification is performed in the descriptor space.

Sample results of our environment reconstruction are shown in Fig. 9. Model texturing is based on a set of photographs taken in the scene.

3. Creating 4D models of walking pedestrians

Relatively small objects such as pedestrians cannot be reconstructed from the LIDAR range data in sufficient detail since the data is too sparse and, in addition, it only gives 2.5D information. Therefore we create properly detailed, textured dynamic models indoors, in a 4D reconstruction studio. The hardware and software components of such a studio can be found in^{4,8}. For completeness, we give below a brief description of the reconstruction process.

Fig. 7 shows a sketch and a panorama of the studio where green curtains and carpet form homogeneous background to facilitate segmentation of the actor. The frame carries 12 calibrated and synchronised video cameras placed uniformly around the scene, and one additional camera on the top in the middle. The cameras are surrounded by programmable LEDs that provide direct illumination. The studio has ambient illumination, as well. Seven PC-s provide the computing power and control the cameras and the lighting.

Currently, each set of 13 simultaneous video frames captured by the cameras is processed independently from the previous one. For a set of 13 images, the system creates a textured 3D model showing a phase of actor’s motion. The main steps of the completely automatic 3D reconstruction process are as follows:

1. Colour images are extracted from the captured raw data.
2. Each colour image is segmented to foreground and background. The foreground is post-processed to remove shadows⁴.
3. A volumetric model is created using the Visual Hull algorithm¹³.
4. A triangulated mesh is obtained from the volumetric model using the Marching Cubes algorithm¹⁶.
5. Texture is added to the triangulated mesh based on triangle visibility⁸.

Fig. 8 shows an example of augmented reality created with the help of the 4D reconstruction studio. Several consecutive phases of an avatar walking in a virtual environment are displayed.

4. Integrating and visualising the spatio-temporal scene model

The last step of the workflow is the integration of the system components and visualisation of the integrated model. The walking pedestrian models are placed into the reconstructed environment so that the center point of the feet follows the trajectory extracted from the LIDAR point cloud sequence. Currently, we use the assumptions that the pedestrians walk

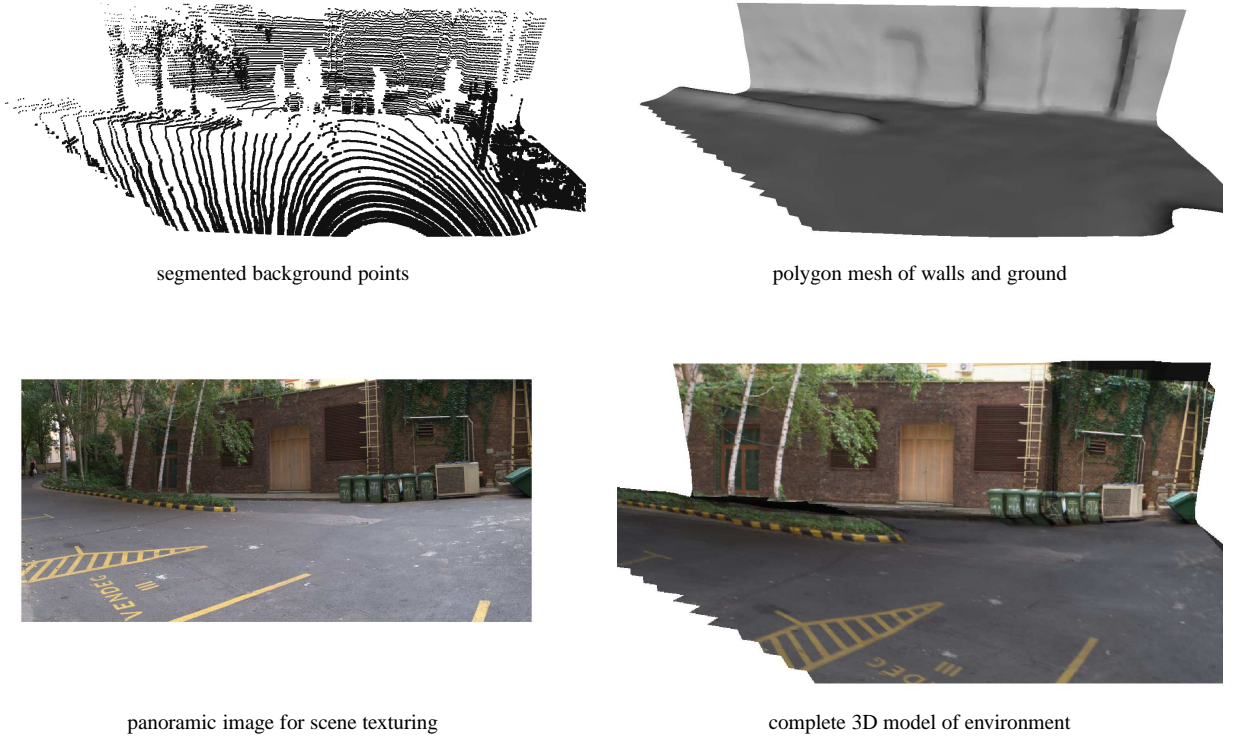


Figure 6: Point cloud segmentation and environment reconstruction.



Figure 7: Sketch and panorama of a 4D reconstruction studio.

forward along their trajectories. The top view orientation of a person is calculated from the variation of the 2D track.

To combine the 3D-4D data of different types arriving in different formats and visualise them in a unified format, we have developed a customised software system. All models are converted to the general-purpose OBJ format²⁰ which is supported by most 3D modelling programs and enables user to specify both geometry and texture.

Our visualisation program is based on the VTK Visualisation Kit¹¹. Its primary goal is to efficiently support combin-

ing static and dynamic models allowing their multiplication and optimising the usage of computational resources. One can easily create mass scenes that can be viewed from arbitrary viewpoint, rotated and edited. Any user interaction with the models, such as shifting and scaling, is allowed and easy to perform.

The dynamic shapes can be multiplied not only in space, but in time, as well. Our 4D studio is relatively small. Typically, only two steps of a walking sequence can be recorded and reconstructed. This short sequence can be multiplied and



Figure 8: A 4D studio actor walking in virtual environment.

seamlessly extended in time to create an impression of a walking person. To achieve this, the system helps the user by shifting the phases of motion in space and time while appropriately matching the sequence of the models.

An important requirement was to visualise the 3D motions of the avatars according to the trajectories provided by the LIDAR pedestrian tracking unit. An avatar follows the assigned 3D path, while rotation of the model to the left or right in the proper direction is automatically determined from the trajectory. Sample final results of the complete 4D reconstruction and visualisation process are demonstrated in Fig. 9.

5. Conclusion and outlook

In this paper, we have introduced a complex system on the interpretation and 4D visualisation of dynamic outdoor scenarios containing multiple walking pedestrians. As a key novelty, we have connected two different modalities of perception: a LIDAR point cloud stream from a large outdoor environment, and an indoor 4D reconstruction studio, which is able to provide detailed models of moving people. The proposed approach points towards real-time free-viewpoint and scalable visualisation of large scenes, which will be a crucial point in future augmented reality and multi modal communication applications. As future plans, we aim to extend the investigations to point cloud sequences collected from a moving platform, and also implement automatic field object recognition and surface texturing modules.

References

1. C. Benedek, Z. Jankó, C. Horváth, D. Molnár, D. Chetverikov, and T. Szirányi. An integrated 4D vision and visualisation system. In *International Conference on Computer Vision Systems (ICVS)*, volume 7963 of *Lecture Notes in Computer Science*, pages 21–30. St. Petersburg, Russia, 2013.
2. C. Benedek, D. Molnár, and T. Szirányi. A dynamic MRF model for foreground detection on range data sequences of rotating multi-beam lidar. In *International Workshop on Depth Image Analysis, LNCS*, Tsukuba City, Japan, 2012.
3. F. Bernardini and et al. Mittleman, J. The Ball-Pivoting algorithm for surface reconstruction. *IEEE Transactions on Visualization and Computer Graphics*, 5(4):349–359, 1999.
4. C. Blajovici, D. Chetverikov, and Zs. Jankó. 4D studio for future internet: Improving foreground-background segmentation. In *IEEE International Conference on Cognitive Infocommunications (CogInfoCom)*, pages 559–564. IEEE, 2012.
5. Q. Chen. Airborne LIDAR data processing and information extraction. *Photogrammetric engineering and remote sensing*, 73(2):109, 2007.
6. R.O. Duda and P.E. Hart. Use of the hough transformation to detect lines and curves in pictures. In *Comm. of the ACM*, volume 15, pages 11–15, 1972.
7. M.A. Fischler and R.C. Bolles. Random sample consensus: A paradigm for model fitting with applications to image analysis and automated cartography. In *Comm. of the ACM*, volume 24, pages 381–395, 1981.
8. J. Hapák, Z. Jankó, and D. Chetverikov. Real-time 4D reconstruction of human motion. In *Proc. 7th International Conference on Articulated Motion and Deformable Objects (AMDO 2012)*, volume 7378 of *Springer LNCS*, pages 250–259, 2012.
9. X. Hu, C.V. Tao, and Y. Hu. Automatic road extraction from dense urban area by integrated processing of high resolution imagery and LIDAR data. *International Archives of Photogrammetry, Remote Sensing and Spatial Information Sciences*, 35:B3, 2004.
10. Z. Jankó, D. Chetverikov, and J. Hapák. 4D reconstruction studio: Creating dynamic 3D models of moving actors. In *Proc. Sixth Hungarian Conference on Computer Graphics and Geometry*, pages 1–7, 2012.
11. Kitware. VTK Visualization Toolkit. <http://www.vtk.org>, 2013.
12. K. Lai and D. Fox. Object recognition in 3D point clouds using web data and domain adaptation. *International Journal of Robotic Research*, 29(8):1019–1037, 2010.

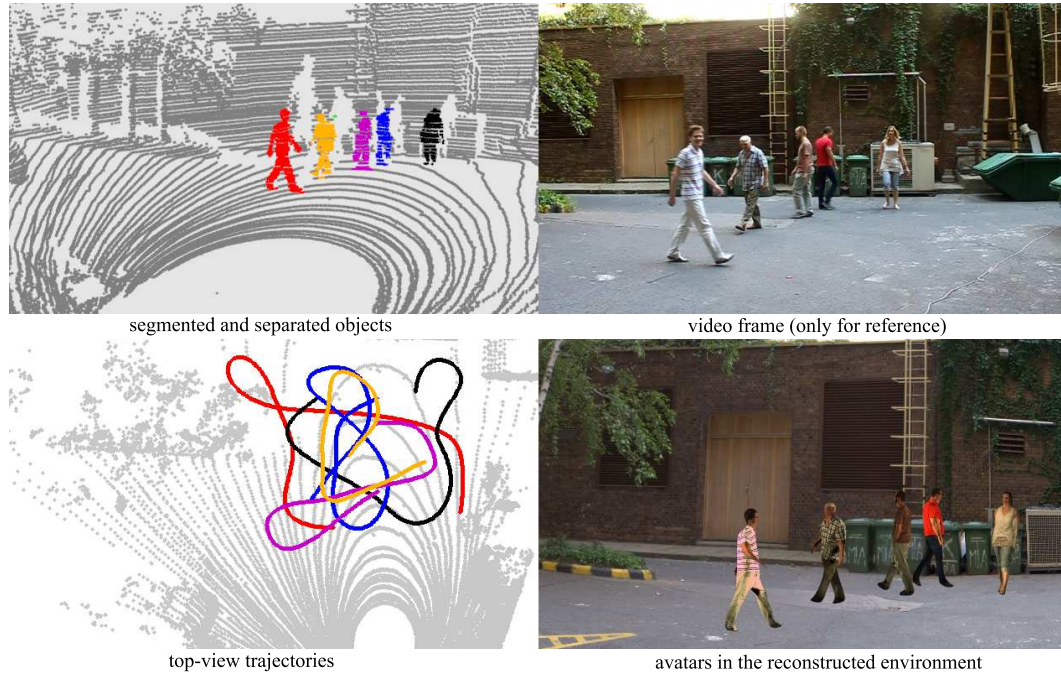


Figure 9: Results of object tracking and integrated dynamic scene reconstruction.

13. A. Laurentini. The visual hull concept for silhouette-based image understanding. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 16:150–162, 1994.
14. H.S. Lee and N.H. Younan. DTM extraction of LiDAR returns via adaptive processing. *Geoscience and Remote Sensing, IEEE Transactions on*, 41(9):2063–2069, 2003.
15. P. Lindner and G. Wanielik. 3D LIDAR processing for vehicle safety and environment recognition. In *IEEE Workshop on Computational Intelligence in Vehicles and Vehicular Systems*, pages 66–71. IEEE, 2009.
16. W.E. Lorensen and H.E. Cline. Marching cubes: A high resolution 3D surface construction algorithm. In *Proc. ACM SIGGRAPH*, volume 21, pages 163–169, 1987.
17. S.C. Popescu and R.H. Wynne. Seeing the trees in the forest: Using LIDAR and multispectral data fusion with local filtering and variable window size for estimating tree height. *Photogrammetric Engineering and Remote Sensing*, 70(5):589–604, 2004.
18. C. Stauffer and W. E. L. Grimson. Learning patterns of activity using real-time tracking. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 22:747–757, 2000.
19. F. Tarsha-Kurdi, T. Landes, and et al. Grussenmeyer, P. Hough-transform and extended RANSAC algorithms for automatic detection of 3D building roof planes from LIDAR data. *International Archives of Photogrammetry, Remote Sensing and Spatial Information Systems*, 36:407–412, 2007.
20. Wavefront Technologies. OBJ file format. Wikipedia, Wavefront .obj file, 2013.