

Video-Based Camera Localization Using Anchor View Detection and Recursive 3D Reconstruction

Hajime Taira^{1,†} Koki Onbe^{1,†} Naoyuki Miyashita^{2,‡} Masatoshi Okutomi^{1,†}

¹Tokyo Institute of Technology ²Olympus Corporation

[†]{htaira,konbe,mxo}@ok.sc.e.titech.ac.jp [‡]naoyuki.miyashita@olympus.com

Abstract

In this paper we introduce a new camera localization strategy designed for image sequences captured in challenging industrial situations such as industrial parts inspection. To deal with peculiar appearances that hurt standard 3D reconstruction pipeline, we exploit pre-knowledge of the scene by selecting key frames in the sequence (called as anchors) which are roughly connected to a certain location. Our method then seek the location of each frame in time-order, while recursively updating an augmented 3D model which can provide current camera location and surrounding 3D structure. In an experiment on a practical industrial situation, our method can localize over 99% frames in the input sequence, whereas standard localization methods fail to reconstruct a complete camera trajectory.

1 Introduction

Determining a location of the camera is one of the fundamental task in computer vision, supporting a growing need of 3D reconstruction such as Structure from Motion (SfM) and Simultaneous Localization and Mapping (SLAM) [1–4]. They are also directly applicable for a video-based navigation system that suggests the temporal user location from a sequential image series.

Such navigation is especially beneficial for some industrial robotics scenarios, where an operator often cannot directly observe the scene [5, 6]. In specific, for an industrial parts inspection, a thin diameter probe (industrial borescope) is inserted to the inner of a product and inspects its damages or defects. Whereas other auxiliary sensors [7, 8] often are not available for practical borescopes, a pure vision-based localization can still be helpful to guess defective locations while associating to their appearances from an image sensor [6]. However, common image-based camera localization techniques [2, 9, 10], which simultaneously estimate the camera location and surrounding 3D structure, often fail to reconstruct a valid model due to peculiar appearances in industrial situations.

In this paper, we attempt to handle such challenging industrial parts inspection scenarios. First, to deal with special appearance in the industrial scene, we employ the 3D structure-based approach [9–14] for localization. Instead of reconstructing cameras for all re-

lated images at once, we efficiently localize input images by preparing a pre-constructed 3D model presenting the targeted scene and registering new images to the model. Second, to stably localize all of video frames captured during inspection, we design our system to register cameras in time-order while incrementally updating 3D model, which makes it easier to find locations of consecutive frames. Also, we employ a new technique based on a typical key-frame (called as *anchor*) in the input sequence that is connected to a certain object in the target scene and contributes to a stable image registration. We finally test our system in one specific inspection scenario for an industrial product and validates its performance in the challenging situation.

Related works. For the input of image series, SfM [1, 2] and SLAM [3, 4, 15, 16] are the well known techniques to reconstruct their 6-dimensional camera poses together with surrounding 3D structure. Whereas most SLAM methods assume a sequential images input and obtain a camera trajectory in time-order, SfM in a recursive manner [17–20] has also been developed to provide a temporal 3D reconstruction in a real-time processing. Since they originally obtain a scale-indeterminate 3D model, several works also estimate cameras in a real-scale by using auxiliary sensors [7, 8, 21–23] or adjusting the model with a known property of the scene [6, 9, 10]. When a pre-constructed 3D model of the targeted scene is also available, camera poses can be further accurately estimated by registering images directly to the model [9–13]. Kroeger *et al.* [24] proposed a video registration scheme to an SfM model. Our method also registers image sequences to an SfM model while incorporating a known property of the scene to deal with challenging industrial scenario.

2 Video frame localization via anchor view detection and recursive 3D reconstruction

Fig. 1 shows the proposed camera localization pipeline for sequential images. We assume a pre-mapped SfM model consisting of 3D scene points and pre-captured (database) cameras as the reference information for the targeted scene (reference model). For the input of an image sequence consisting of consecutive image frames, we first find a key frame connected to a specific location

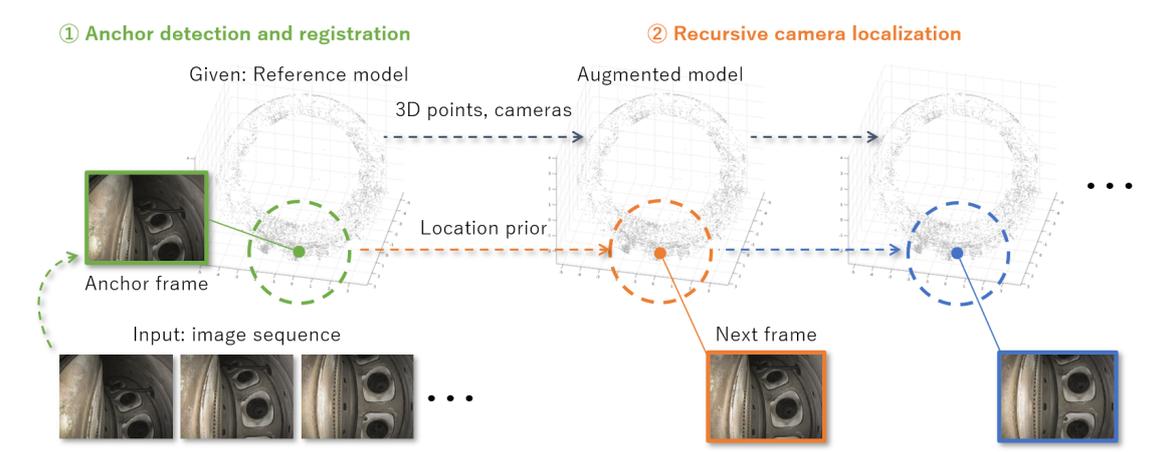


Figure 1. The whole camera localization pipeline for a sequential image series. We start to reconstruct cameras from one specific anchor frame capturing the characteristic location in the targeted scene. Then we localize remaining frames by incrementally updating an augmented model while registering images neighbor to the previously reconstructed camera.

in the reference model (Sec. 2.1). Next we sequentially localize remaining images via an SfM-like reconstruction scheme (Sec. 2.2) which recursively update a 3D model by registering new images. The reconstruction will continue until all input images have been registered.

2.1 Anchor-based camera localization

The major failure cases of reconstruction in industrial scenarios are often be attributed to the failure of localization due to the lack of texture, or highly frequent objects such as standardized industrial parts. To deal with these challenging appearance, we attempt to first register some key frames (anchors) which locations can be roughly determined by the pre-knowledge of the targeted scene, *e.g.*, capturing an unique object or a marker in the scene. In the later section (Sec. 3) we will describe our CNN-based anchor detector trained for one specific situation of industrial parts inspection, which constructs a subset of images potentially to be anchors out of the input sequence. Please note that in more general cases, such anchor frames can either be specified manually, or automatically detected via any object recognizer. The detected anchors can be relatively easily registered to the reference model, and also can be used as a spatial guide of other frames. As described in the later section (Sec. 2.2), this anchor-based approach is particularly beneficial to stably localize sequential images.

Anchor registration to the reference model. We register the detected anchor images to the reference model via a standard SfM scheme constructing a temporal augmented model including anchors. We first seek a subset of database images that share views with

anchors, and perform pairwise local feature matching towards each of anchor images [25]. Consequent PnP and point triangulation steps [2] obtain the initial camera pose of anchors and 3D scene points corresponds to the local features in these newly registered images. The camera poses of anchors are finally refined by the bundle adjustment [2] which minimizes the reprojection errors of 3D scene points with respect to their local observations in the images. Please note that in this step we freeze parameters of database cameras and existing scene points in the reference model, so this refinement does not affect the consistency of the augmented model and referenced SfM model presenting the map information of the scene.

2.2 Recursive 3D reconstruction for sequential images

After registering anchor frames to the SfM model, we start to determine the location of all images by incrementally register each of frames, while updating the augmented model.

Feature matching towards spatial and temporal neighbors. The reconstruction begins from the consecutive frame of the earliest anchor frame, and continue in the time-stamp order. First we perform a *spatially-guided* feature matching that seek correspondences between local observations in the current frame and cameras in the current augmented model. Instead of the location knowledge used in Sec. 2.1, this time we exploit the most recently registered camera of the input sequence as a more precise location prior. As in anchor registration, we collect a set of database images that are spatially neighbor to the recent camera and perform pairwise local feature matching towards the cur-

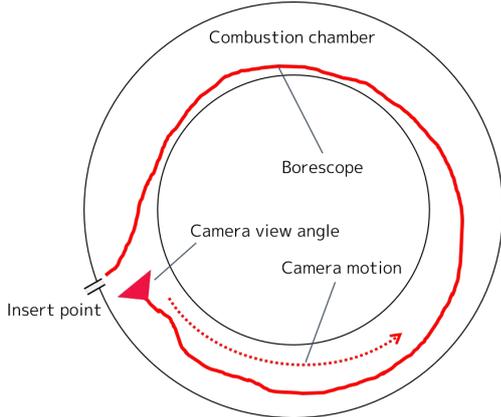


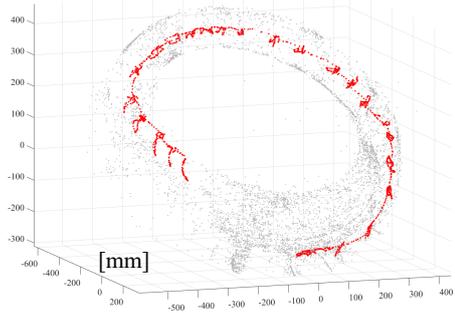
Figure 2. A rough sketch of the inspection for a combustion chamber.

rent frame. Additionally, we select 20 best database images via image retrieval [26] and also match them towards the current frame.

To increase the number of local connectivity of the current frame towards the augmented model, we also perform a *temporally-guided* feature matching that matches the current frame towards the neighbor frames, which have highly relevant appearances thus are easy to match. We make a set of N neighbor frames that have previously been registered to the model and match them towards the current frame. In our experiments, we set N as 25 frames, which can remain a precise reconstruction with a controllable cost, but please be sure that this setting should be determined for each specific scene while also considering the video frequency.

Recursive camera localization for image sequence. From the set of local feature matches, we can also extract correspondences between local features and the existing 3D scene points. If the sufficient number of these 2D-to-3D matches exists, the 6D camera pose of the current frame can be obtained by solving the Perspective-n-Points problem [27,28]. Then we triangulate new 3D scene points seen from the current frame, which can also help to register the next frame.

We finally perform a bundle adjustment [29] to refine the newly added scene points and cameras. As in anchor registration, we freeze 3D points and cameras existing in the reference model, and refine cameras for all input frames and new 3D points while minimizing the reprojection errors of the model. For efficiency, we perform this refinement every 10 new frames have been registered. After all, we get a new temporal SfM model that provides camera location of the recent input frames, which is used to register next consecutive frames. The system continues to register new frames until the end of the input video, and finally obtains an SfM model including 6D cameras of all input frames.



Anchor



Figure 3. **The test environment and image samples.** Top: Gray dots represent the 3D reference model of the combustion chamber in a jet engine, while red dots are the ground-truth camera location of the test sequence. Bottom: Video frames in the test sequence.

Table 1. Localization results.

Method	#Cameras	MAE	Median error
		[mm]	[mm]
Single image	1,450 (63.6%)	1.34	0.60
Ours	2,247 (98.6%)	2.44	0.42

3 Experiments

Test scene: The inside of a jet engine. We test our localization system during the inspection for a combustion chamber in an aircraft jet engine [30]. The inspection is usually done by inserting an industrial borescope through an insert point as to go around the chamber (*c.f.* Fig. 2). Then the inspector observes the inside of the chamber via a monocular camera of the borescope, while pulling out the borescope. We capture three image sequences during independent trials, two as the database sequences for constructing the initial SfM model, and the other as the query sequences for testing. The reference model is constructed by COLMAP [2], resulting in a model consisting of 5,107 cameras and 1.5M scene points. Separately, we gather the ground-truth cameras of query sequences by constructing another SfM model for the query and one of the database sequence. The reconstructed query cameras are registered to the reference model by estimating the similarity transform between the shared frames [31], resulting in 2,254 query frames annotated with the ground-truth location, out of 2,280 query frames (*c.f.* Fig. 3).

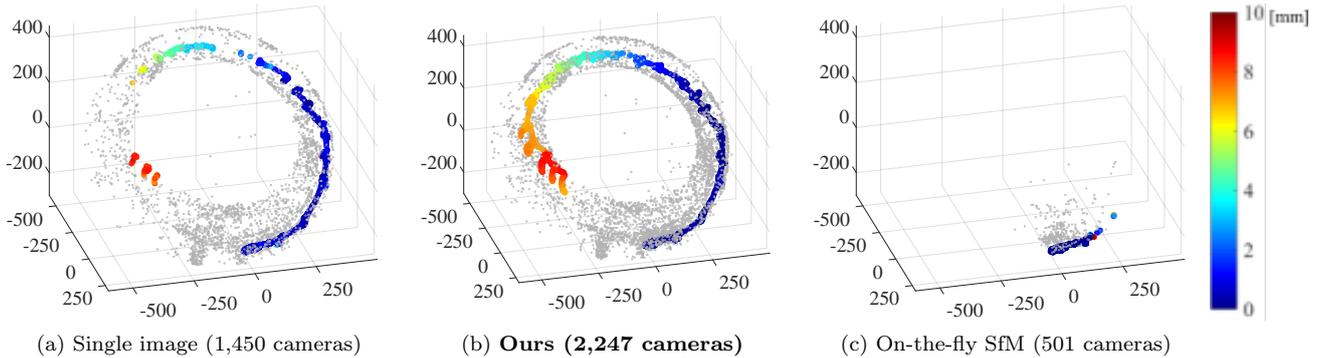


Figure 4. **Reconstructed camera trajectory.** Each camera is colored by its positional error (mm).

CNN anchor detector. According to the trajectory of the borescope illustrated in Fig. 2, the beginning of video often captures the insert point (and the inserted borescope itself), which has a special appearance compared to other part of the chamber (*c.f.* Fig. 3). Therefore we build a CNN classifier to detect such unique frames as anchors. The model is based on ResNet18 [32] architecture, while discarding last two layers (conv4 and conv5, namely) and modifying the final fully connected layer to obtain an one dimensional score of the input image to be an anchor. We manually annotate 518 database images (roughly 10% of all frames) which see the insert point and train the model by minimizing a standard margin loss [33]. In the testing phase, we feed each of input frames to the trained model, and gather a subset of 32 anchor frames.

Implementation. We implement our method mainly based on COLMAP [2], a well-known SfM tool. We modify bundle adjustment [29] so that we can freeze the reference model. The anchor detector is implemented using PyTorch library.

Results. As the main comparison opponent to our method, we also evaluate a baseline localization method for single image: For each input frame, we gather 20 similar database images via image retrieval and match their local features. Using these correspondences, the camera pose of each frame is independently estimated by solving PnP.

Tab. 1 reports the statistics of the localized cameras. Our method can localize further more cameras than the single image method, while remaining the accuracy. Fig. 4 also plots the localized cameras in the model. Due to the challenging scene nature and partly incomplete reconstruction of the reference model, single image method (a) fails to find locations of the latter part of the input sequence. On the other hand, our method (b) achieves a continuous camera trajectory within acceptable errors, which prove the dominance of our approach using multiple sequential images to support each other location.

As an alternative approach which performs *on-the-fly* scene reconstruction [9, 10], we also construct an

SfM model via standard incremental SfM implemented by COLMAP [2], using only input frames. The model is then registered to the scene by estimating the similarity transform between reconstructed cameras and their ground-truth location. This approach, however, reconstructs only part of the scene seen from few images (Fig. 4 (c)). This result clearly points the fact that the well known camera tracking approach that simultaneously estimates 3D structure and camera locations, including “popular” SLAM approaches, actually cannot deal with severe appearances such as in industrial scenes, whereas our method can still obtain an accurate camera trajectory while exploiting the pre-constructed 3D model.

4 Conclusion

In this paper, we have proposed a new camera localization system designed for an image sequence captured in the challenging industrial scene. Our method starts reconstruction from a reliable anchor frame that captures an unique object in the scene, and sequentially register neighbor frames while exploiting recently registered frames as the location prior for the new frame. In the experiment on an industrial parts inspection scenario, the proposed method achieves an accurate and stable camera trajectory whereas other methods can localize only a part of the sequence. We believe our recursive 3D reconstruction anchored to any static object in the scene, is also beneficial for localization problems in many robotics situations, where the operator can perceive the current location and the surrounding environment through the augmented 3D model obtained during the reconstruction. One of the future work would be to achieve a real-time processing for a sequential image stream from a camera, so that the system can provide the augmented model in a practical timing.

Acknowledgement. This work is partly supported by JSPS KAKENHI Grant Number 17H00744.

References

- [1] C. Wu, “Towards Linear-time Incremental Structure From Motion,” in *Proc. 3DV*, 2013.
- [2] J. L. Schönberger and J.-M. Frahm, “Structure-From-Motion Revisited,” in *Proc. CVPR*, 2016.
- [3] C. Cadena, L. Carlone, H. Carrillo, Y. Latif, D. Scaramuzza, J. Neira, I. Reid, and J. J. Leonard, “Past, present, and future of simultaneous localization and mapping: Toward the robust-perception age,” *IEEE Trans. Robotics*, vol. 32, no. 6, pp. 1309–1332, 2016.
- [4] R. Mur-Artal, J. M. M. Montiel, and J. D. Tardos, “ORB-SLAM: A versatile and accurate monocular SLAM system,” *IEEE Trans. Robotics*, vol. 31, no. 5, pp. 1147–1163, 2015.
- [5] P. Hansen, H. Alismail, B. Browning, and P. Rander, “Stereo Visual Odometry for Pipe Mapping,” in *Proc. IEEE/RSJ Conf. on Intelligent Robots and Systems*, 2011.
- [6] S. Kagami, H. Taira, N. Miyashita, A. Torii, and M. Okutomi, “3D Pipe Network Reconstruction Based on Structure from Motion with Incremental Conic Shape Detection and Cylindrical Constraint,” in *Proc. ISIE*, 2020.
- [7] S. Esquivel, R. Koch, and H. Rehse, “Reconstruction of sewer shaft profiles from fisheye-lens camera images,” in *Joint Pattern Recognition Symposium*, 2009.
- [8] P. Hansen, H. Alismail, P. Rander, and B. Browning, “Visual mapping for natural gas pipe inspection,” *Intl. J. of Robotics Research*, vol. 34, no. 4-5, pp. 532–558, 2015.
- [9] T. Sattler, A. Torii, J. Sivic, M. Pollefeys, H. Taira, M. Okutomi, and T. Pajdla, “Are large-scale 3D models really necessary for accurate visual localization?” in *Proc. CVPR*, 2017.
- [10] A. Torii, H. Taira, J. Sivic, M. Pollefeys, M. Okutomi, T. Pajdla, and T. Sattler, “Are Large-Scale 3D Models Really Necessary for Accurate Visual Localization?” *IEEE PAMI*, vol. 43, no. 3, pp. 814–829, 2021.
- [11] T. Sattler, B. Leibe, and L. Kobbelt, “Efficient & Effective Prioritized Matching for Large-Scale Image-Based Localization,” *IEEE PAMI*, vol. 39, no. 9, pp. 1744–1756, 2017.
- [12] H. Taira, M. Okutomi, T. Sattler, M. Cimpoi, M. Pollefeys, J. Sivic, T. Pajdla, and A. Torii, “InLoc: Indoor visual localization with dense matching and view synthesis,” in *Proc. CVPR*, 2018.
- [13] —, “InLoc: Indoor visual localization with dense matching and view synthesis,” *IEEE PAMI*, vol. 43, no. 4, pp. 1293–1307, 2021.
- [14] C. Toft, W. Maddern, A. Torii, L. Hammarstrand, E. Stenborg, D. Safari, M. Okutomi, M. Pollefeys, J. Sivic, T. Pajdla, F. Kahl, and T. Sattler, “Long-term visual localization revisited,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020.
- [15] J. Engel, T. Schöps, and D. Cremers, “LSD-SLAM: Large-Scale Direct Monocular SLAM,” in *Proc. ECCV*, 2014.
- [16] S. Sumikura, M. Shibuya, and K. Sakurada, “Openvslam: A versatile visual slam framework,” in *Proc. ACM Multimedia*, 2019, pp. 2292–2295.
- [17] F. Nyberg and A. Heyden, “Recursive Structure from Motion Using Hybrid Matching Constraints with Error Feedback,” in *Dynamical Vision*. Springer, 2006, pp. 285–298.
- [18] S. Bronte, M. Paladini, L. M. Bergasa, L. Agapito, and R. Arroyo, “Real-Time Sequential Model-Based Non-Rigid SFM,” in *Proc. IEEE/RSJ Conf. on Intelligent Robots and Systems*, 2014.
- [19] S. Song and M. Chandraker, “Robust Scale Estimation in Real-Time Monocular SFM for Autonomous Driving,” in *Proc. CVPR*, 2014.
- [20] B. Resch, H. Lensch, O. Wang, M. Pollefeys, and A. Sorkine-Hornung, “Scalable Structure from Motion for Densely Sampled Videos,” in *Proc. CVPR*, 2015.
- [21] Z. Zhang, R. Zhao, E. Liu, K. Yan, and Y. Ma, “Scale Estimation and Correction of the Monocular Simultaneous Localization and Mapping (SLAM) Based on Fusion of 1D Laser Range Finder and Vision Data,” *Sensors*, vol. 18, no. 6, p. 1948, 2018.
- [22] S. Sumikura, K. Sakurada, N. Kawaguchi, and R. Nakamura, “Scale Estimation of Monocular SfM for a Multimodal Stereo Camera,” in *Proc. ACCV*, 2018.
- [23] Z. Jiang, H. Taira, N. Miyashita, and M. Okutomi, “VIO-Aided Structure from Motion Under Challenging Environments,” in *Proc. ICIT*, 2021.
- [24] T. Kroeger and L. Van Gool, “Video Registration to SfM Models,” in *Proc. ECCV*, 2014.
- [25] D. G. Lowe, “Distinctive image features from scale-invariant keypoints,” *IJCV*, vol. 60, no. 2, pp. 91–110, 2004.
- [26] J. L. Schönberger, T. Price, T. Sattler, J.-M. Frahm, and M. Pollefeys, “A Vote-and-Verify Strategy for Fast Spatial Verification in Image Retrieval,” in *Proc. ACCV*, 2016.
- [27] M. A. Fischler and R. C. Bolles, “Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography,” *Comm. ACM*, vol. 24, no. 6, pp. 381–395, 1981.
- [28] L. Quan and Z. Lan, “Linear n-point camera pose determination,” *IEEE PAMI*, vol. 21, no. 8, pp. 774–780, 1999.
- [29] S. Agarwal, K. Mierle, and Others, “Ceres solver,” <http://ceres-solver.org>.
- [30] “CFM56,” www.cfm aeroengines.com/engines/cfm56/.
- [31] S. Umeyama, “Least-squares estimation of transformation parameters between two point patterns,” *IEEE PAMI*, no. 4, pp. 376–380, 1991.
- [32] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proc. CVPR*, 2016.
- [33] R. Arandjelović, P. Gronat, A. Torii, T. Pajdla, and J. Sivic, “NetVLAD: CNN Architecture for Weakly Supervised Place Recognition,” *IEEE PAMI*, vol. 40, no. 6, pp. 1437–1451, 2018.