# TopoTag: A Robust and Scalable Topological Fiducial Marker System

Guoxing Yu      Yongtao Hu      Jingwen Dai

Guangdong Virtual Reality Technology Co., Ltd.

{calvin.yu, ythu, dai}@ximmerse.com

## Abstract

*Fiducial markers have been playing an important role in augmented reality (AR), robot navigation, and general applications where the relative pose between a camera and an object is required. We introduce TopoTag, a robust and scalable topological fiducial marker system, which supports reliable and accurate pose estimation from a single image. TopoTag uses topological and geometrical information in marker detection to achieve higher robustness. Without sacrificing bits for higher recall and precision like previous systems, TopoTag can use full bits for ID encoding and supports tens of thousands unique IDs and easily extends to millions and more by adding more bits, thus achieves perfect scalability. We collect a large dataset including in total 169,713 images for evaluation, involving in-plane and out-of-plane rotation, image blur, different distances and various backgrounds, etc. Experiments show that TopoTag significantly outperforms previous fiducial marker systems in terms of various metrics, including detection accuracy, vertex jitter, pose jitter and accuracy, etc. In addition, TopoTag supports occlusion as long as main tag topological structure is maintained and flexible shape design where users can customize inter and outer marker shapes. Our dataset, marker design and detection algorithm are public to the community[1].*

## 1. Introduction

In this paper, we introduce TopoTag, a new fiducial marker and detection algorithm that is more robust and accurate than current fiducial marker systems. Fiducial markers are artificial objects designed to be easily detected in an image from a variety of perspectives. They are widely used for augmented reality and robotics applications because they enable localization and landmark detection in featureless environments [16]. Previous works on fiducial markers mainly focus on one or more of the following areas: (1) improve detection accuracy via specialized tag de-

---

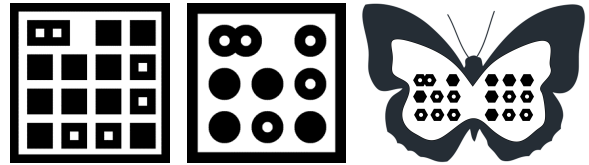[1]TopoTag homepage: https://sites.google.com/view/topotag/home



Figure 1: Three TopoTag markers.

sign [17, 18, 1, 32]; (2) decrease pose estimation error via precise vertex estimation [45] or introducing more feature points [5]; (3) increase unique identities [9, 19, 15, 21]; (4) improve robustness under occlusion [5, 20] and other use cases [35, 8, 14] and (5) speedup [30, 45, 16, 39].

TopoTag utilizes topological information in tag design to improve robustness, which achieves perfect detection accuracy on the large dataset we collected and also datasets from others. We show that all tag bits can be used to encode identities without sacrificing detection accuracy, thus achieving rich identification and easy scalability for more. In addition, TopoTag offers more feature correspondences for better pose estimation. Results show that TopoTag achieves the best performance in vertex jitter, pose error and jitter. TopoTag also supports occlusion as long as main tag topological structure is maintained and flexible shape design where users can customize inter and outer marker shapes.

We collect a large dataset including in total 169,713 images with TopoTag and several state-of-the-art tag systems. Robot arm is used to guarantee repeated trajectories for each tag for fair comparison. The rich modalities of the dataset including in-plane and out-of-plane rotations, image blur, different distances and various backgrounds, etc offer a challenging benchmark to evaluate.

In summary, the contributions of this paper are: (1) we present TopoTag, a topological-based fiducial marker and detection algorithm; (2) we demonstrate that our tag achieves the best performance in various metrics including detection accuracy, localization jitter and accuracy, etc and support occlusion and flexible shapes; (3) we show that it's possible in tag design to use full bits for ID encoding without sacrificing detection accuracy, thus achieving easy scal-
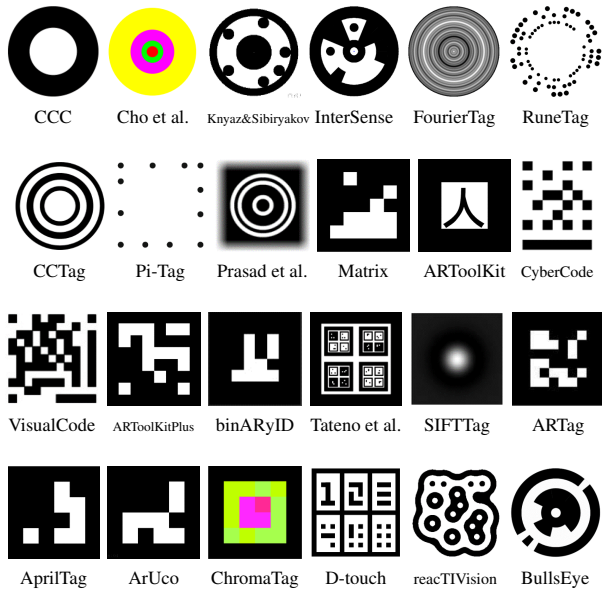
Figure 2: Existing fiducial marker systems.

ability; and (4) we collect a large dataset of various tags, involving in-plane and out-of-plane rotation, image blur, different distances and various backgrounds, etc.

## 2. Related Work

Figure 2 shows many different fiducial marker systems discussed in this section.

**Circular patterns.** Among the earliest work, Gatrell et al. [22] propose to use concentric contrasting circle (CCC). It's further enhanced in [9] by adding colors and multiple scales. In [27, 31], dedicated data rings are added to the marker design for rich identification. Sattar et al. [41] and Xu et al. [46] propose FourierTag with a frequency image as the signature. RuneTag [5, 4] and Pi-Tag [6] propose using rings of dots to improve robustness to occlusion and provide more points for pose estimation. CCTag [7, 8] and followed work by Prasad et al. [35] use multiple rings to increase robustness to blur and ring width for encoding.

**Square patterns.** To be easily localized, the majority of fiducial systems are designed to contain a thick square border. Matrix [36], CyberCode [37] and VisualCode [38] are the first and simplest proposals. ARToolkit [25] is well known and widely used in many virtual reality applications. It includes a pattern in their inner region for identification via image correlation. ARTag [17] and ARToolkitPlus [44] improve the recognition technique with a binary coded pattern. In addition, they are designed with error correction mechanism to increase robustness. BinARyID [19] proposes a method to generate markers focused on avoiding rotation ambiguities. Schweiger et al. [42] propose using

SIFT and SURF filters that are specifically designed for SIFT and SURF detectors. Tateno et al. [43] propose using nested markers to improve performance under different distances. Several works investigate using multiple fiducial markers in a checkerboard to improve camera calibration [1] and reduce the perspective ambiguity by further adding color [15]. AprilTag [32, 45] is a faster and more robust reimplementation of ARTag. Garrido-Jurado et al. [20, 21] propose ArUco using mixed integer programming to generate markers. ChromaTag [16] uses color over AprilTag to improve marker detection speed. Square markers choose to use four corner points for pose estimation, which is the minimum number for unambiguous pose estimation [33]. As a contrast, TopoTag offers more feature correspondences for a better pose estimation.

**Topological patterns.** D-touch [13, 12] is the earliest work to use topological patterns in tag design. Marker detection is based on the region adjacency tree information. D-touch employs a single topology for all markers in the set and does not provide a specific method for computing location and orientation. ReacTIVision [3, 2, 24] improves over D-touch to provide unique identities purely by the topological structure. BullsEye [26] is specially designed optimized for GPU. Both ReacTIVision and BullsEye can only recover 2D location and orientation due to the limitation of not enough matched feature points.

**Machine learning.** Claus et. al [10, 11] use trained classifiers to improve detection in cases of bad illumination and blurring caused by fast camera movement. Randomized forests are also used to learn and detect planar objects [28, 34]. In practice, these algorithms do not achieve detection accuracies on par with detection algorithms specifically designed for marker detection [16].

## 3. TopoTag Design

TopoTag utilizes topological structure information in tag design, which has been proven with robustness with illuminate variation and false detection [13]. Existing fiducial marker systems, especially square patterns, sacrifice tag bits to handle four rotation ambiguities during decoding [15]. Additional bits will also be reserved for incorporating Hamming distance strategy in order to improve false positive rejection. High robustness with topological design helps saving tag bits for encoding identities. To avoid rotation ambiguities, TopoTag introduces *baseline node* in its topological structure. The baseline node is specially designed to be different from other nodes in the tag. TopoTag uses a black node with two white children nodes inside as the baseline node and other black nodes with at most one white child node as *normal nodes*. Note that, baseline node can be defined with other forms like with three or more white children nodes for different needs. Baseline node defines the search starting position of the tag, thus avoiding checking
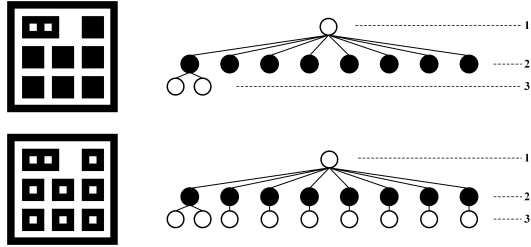
Figure 3: Topological tree of two TopoTags.

rotation ambiguities. All normal nodes are used for identity encoding with 0 denoting no child node and 1 otherwise. The identity encoding for the two markers shown in Figure 3 is 0000000 = 0 and 1111111 = 127 respectively.

For pose estimation, instead of using only four border points in previous square systems [17, 44, 32, 20, 21, 45, 16] which is the minimum number required, TopoTag offers more point correspondences for more accurate pose estimation. Baseline node (more specifically its two children nodes) and all normal nodes are all used as feature points, thus achieving a better pose estimation.

Note that, as TopoTag design is based on topological information, there is no restriction for the shapes used in the tag. Both inner and outer shapes can be customized as long as the desired topological structure is preserved. Figure 1 shows three different design samples of TopoTag. For easy searching and model simplicity, in current TopoTag design, we choose to place all inner nodes uniformly spaced and compacted into a $n \times n$ squared shape.

## 4. TopoTag Detection

Figure 4 outlines main steps of TopoTag detection. Topological information is extensively used for 2D marker detection, and further corresponding geometrical information for ID decoding. 3D pose estimation is achieved by taking advantage of all TopoTag vertices.

### 4.1. 2D Marker Detection

**Threshold map estimation.** Similar to the idea of adaptive thresholding, we estimate threshold for each pixel by analyzing its surrounded pixels. The analysis can be done on the original image. In practice, we find that doing on a downsampled version (scalar $s_1$) is more accurate due to image noise and blur, which also brings speed benefits. Any pixel will be set to $\alpha$ if its value is less than $\alpha$ to remove too dark pixels. Average values are computed on a local region (window size $w$) on the downsampled image. To further handle noise, the downsampled average map can be further downsampled (scalar $s_2$). The final threshold map is achieved by upsampling the downsampled average map by $s_1 \times s_2$ using bilinear interpolation, see Figure 4b.
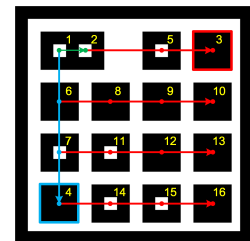
**Binarization.** Binarization is achieved by comparing the input image with the threshold map. A minimum brightness ($\beta$) is set to filter too many small regions (i.e. set to black if pixel value is less than $\beta$). See Figure 4c for an example of binarization result.

**Topological filtering.** After segmentation, we build the topological tree of the connected binary regions. To find candidate tags, we search the tree based on two conditions: (1) the number of children nodes should be within $[\zeta_{\min} - \tau, \zeta_{\max} + \tau]$, where $\zeta_{\min}$ is the number of nodes for tag ID = 0 with all black leaves except the baseline node and $\zeta_{\max}$ for the tag with maximum ID with no black leaves, and $\tau$ is the tolerance level allowed; (2) max depth of the tree should be exactly 3. See Figure 3 for examples of the topological trees for both $\zeta_{\min}$ and $\zeta_{\max}$ cases of our 9-bit tags. Figure 4d shows the result after topological filtering.

**Error correction.** There are possible error nodes within the tag region due to noise or occlusion. Figure 4d shows an example of one error node close to the baseline node because of one ant sitting on the tag. To correct these error nodes, we filter out smaller nodes if their areas are less than $\theta_1\%$ of the baseline node area. Figure 4e shows the result after error correction.

### 4.2. ID Decoding

To decode ID, we need to determine the node sequence and map it to a binary code string. Take a 16-bit TopoTag as an example, see the right figure and Figure 4f of the sequence we find for each node of the tag. To start, we first find the baseline node (including p1 and p2) and determine



its search direction based on whether there are nodes along the direction with angel tolerance $\theta_2$, i.e. p1→p2. Along the direction, we find the node with the largest distance, i.e. p3. For the remaining nodes, we first find the node with largest angle against the baseline direction p1→p2 and then the largest distance along the direction, i.e. p4. p5 is determined along direction p1→p3, p6 and p7 along p1→p4. The remaining nodes are determined in order in the similar way. After finding each node, we can simply map each node to 1 or 0 depending wether it contains a white child node or not and decode the tag based on the binary code string. For the example shown here and in Figure 4, the binary code string is 10000011000110, which is decoded with ID = 8390. It's worthy noting that, ID decoding is processed on the images after removing the perspective distortion to improve robustness.
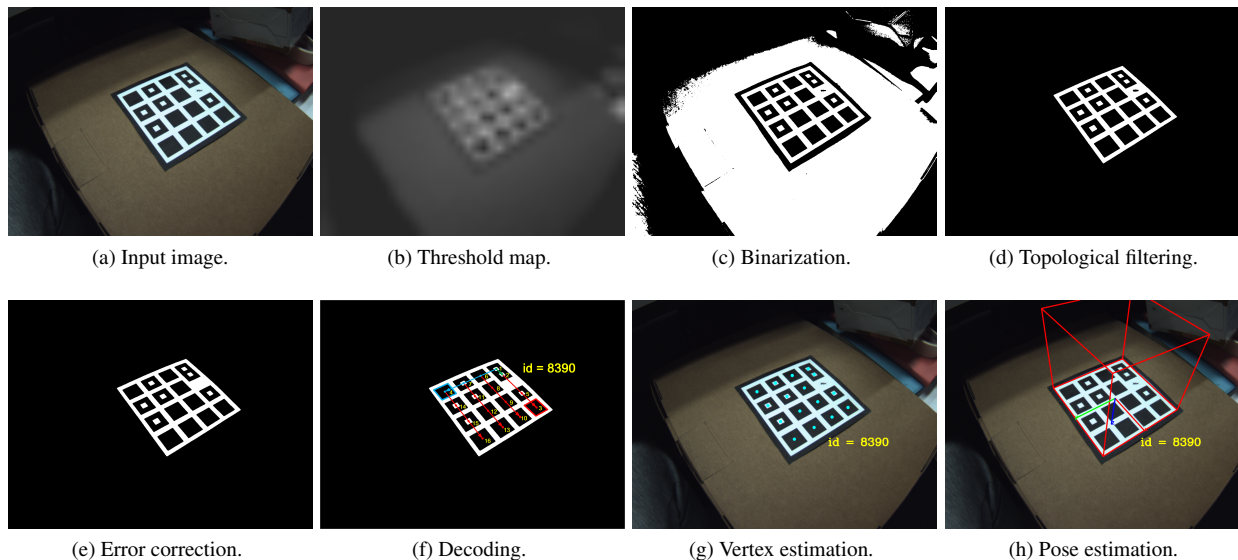
| (a) Input image. | (b) Threshold map. | (c) Binarization. | (d) Topological filtering. |

| (e) Error correction. | (f) Decoding. | (g) Vertex estimation. | (h) Pose estimation. |

Figure 4: Main steps of TopoTag detection. Best viewed in color.

### 4.3. 3D Pose Estimation

For each node, we estimate the vertex by computing the centroid on the original image of its supporting region. The supporting region can be the binary mask or its dilated version (with dilate size $\delta$). The centroid can be determined via image moments, i.e. $\{\bar{u}, \bar{v}\} = \{\frac{M_{10}}{M_{00}}, \frac{M_{01}}{M_{00}}\}$.

For pose estimation, the exact correspondence between the 2D image features and the features of the associated model is needed (feature correspondence). At least 4 points are needed to recover unambiguous pose estimation for planar tags [33]. Unlike most of previous work using only 4 corner points, all TopoTag vertices of tag bits are used for a better pose estimation. For 16-bit tag, 16 vertex correspondences are used, including 2 baseline white nodes and 14 normal black nodes. 6-DoF pose estimation is achieved by solving the PnP problem and Levenberg-Marquardt algorithms [29, 23] based on these feature correspondences.

## 5. Results and Discussion

**Algorithm setup.** Throughout the experiment, we use $s_1 = 4, s_2 = 8, w = 5, \alpha = 45, \beta = 50$ for segmentation, $\tau = 0, \theta_1 = 30, \theta_2 = 0.1$ rad for decoding, $\delta = \max\{2, \lfloor \frac{l}{10} \rfloor\}$ for vertex estimation, where $l$ is the short length of the binary mask region.

All the experiments have been performed on a typical laptop PC equipped with an Intel Core i7-7700HQ processor (8 cores @2.8Ghz) and 8GB of RAM.

### 5.1. Dataset

Previous work, like [32, 5], mainly focus on evaluating performance on synthetic images. Although several work evaluate part of their performance on more realistic scenes, e.g. ARToolKitPlus [44] evaluates the speed on several handheld devices and AprilTag [32, 45] evaluates false positive on LabelMe [40] dataset which is designed for general object detection and recognition research, there is still no uniform dataset for fiducial marker evaluation, thus make it hard to reproduce their result and compare with others. More recently, in ChromaTag [16] work, they collected a dataset to compare their work with AprilTag [32], CCTag [8], and RuneTag [5]. However, different tags are placed side by side during their dataset collection, thus make it not ideal for comparison especially when tags viewed from a large angle as different markers will have different distances and facing angles towards the camera.

In this work, we try to fill this gap by collecting a large dataset, including in total 169,713 images, which involves in-plane and out-of-plane rotations, image blur, various distances and cluttered backgrounds, etc. We use a global shutter industrial camera with 1280×960 resolution streaming at 38.8 fps and 98° diagonal field of view. The exposure is fixed at 10 ms. Using relative long exposure guarantees enough brightness of the captured images, which at the same time introduces image blur phenomenon for more challenging use cases (see the first image in Figure 7 for an example). Camera is fixed to a robot arm[2] to guarantee repeated trajectories for different tags. Figure 5 shows the

---

[2]We use a robot arm from DENSO (VS-6556). Link: https://www.denso-wave.com/en/robot/product/five-six/vs.html
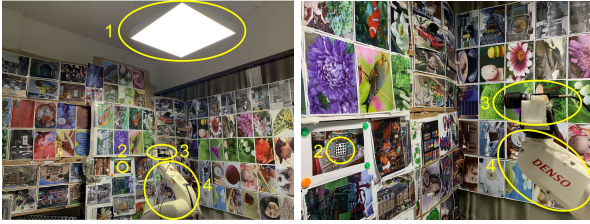
Figure 5: Dataset collection setup. We collect dataset by putting tags (label #2) in a rich textured background within an indoor environment with fixed lighting (label #1). Camera (label #3) is fixed to a robot arm (label #4) to guarantee repeated trajectories for different tags.
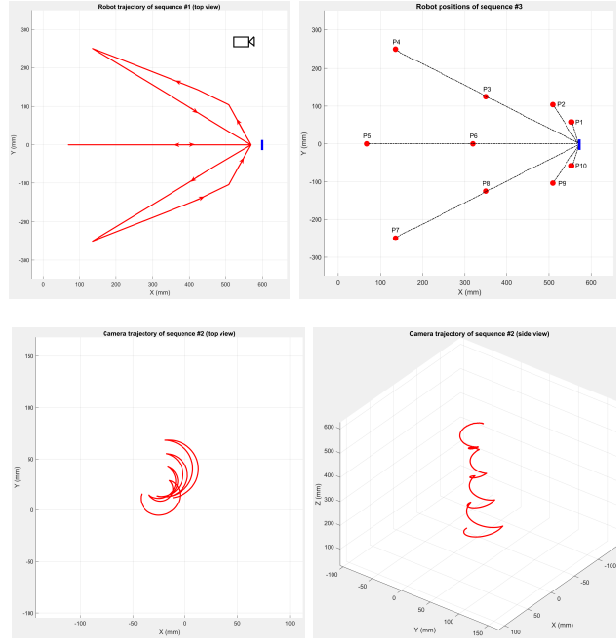


Figure 6: Robot arm trajectory/points in different sequences (top-left for Seq #1, top-right for Seq #3). Camera trajectory is shown for Seq #2 for better visualization (bottom row). Tag position is shown in blue. Best viewed in color.



Figure 7: Sample collected images. Images are from Seq #1 (with ARToolKit), Seq #2 (with AprilTag 25h9) and Seq #3 (with TopoTag) from left to right respectively.

dataset collection setup. Three sequences will be collected for each tag, with trajectory for each sequence shown in Figure 6. In all three sequences, the camera keeps facing the front as the camera icon shown in the first image of Figure 6. In Seq #1, camera moves along several lines at a constant speed, with different out-of-plane rotations for each line including $0°$ (i.e. camera faces the tag in right ahead), $30°$ and $60°$. In Seq #2, the camera moves back and at the same time rotates in-plane within $0\text{-}180°$ at a constant speed back and forth. Note that, as we can only rotate around the robot arm and there is a offset between camera with it, camera's trajectory will not be an ideal half circle. In Seq #3, camera is placed at 10 fixed positions (P1→P10). Besides $0°$, $30°$ and $60°$ out-of-plane rotations as in Seq #1, we further collect data with $75°$ (P1 and P10). In all three sequences, the background is filled with rich textured images to simulate more complex use scenarios.

We collect the dataset for TopoTag and previous tags including ARToolKit [25], ARToolKitPlus [44], ArUco [20], RuneTag [5], ChromaTag [16] and AprilTag [45]. A 16-bit TopoTag is used throughout the experiment as it provides the state-of-the-art most unique identities, see Table 1 for details. And, without loss of generality, the tag comes with square inner and outer shapes, see the first image in Figure 1. For systems with multiple tag families, we collect data for each. For each tag family, we randomly select one ID for evaluation. In our experiment, we randomly select ID = 1 for ARToolKit, 262 for ARToolKitPlus, [104, 90, 136] for ArUco's [16h3, 25h7, 36h12], 107 for RuneTag, 0 for ChromaTag, [0, 204, 25, 1314, 343] for AprilTag's [16h5, 25h7, 25h9, 36h9, 36h11] and 278 for TopoTag. Outer border sizes of all tags are the same of 5 cm. For each tag, there are ≈100,000 images collected, including ≈1,000 for Seq #1, ≈1,200 for Seq #2 and ≈7,800 for Seq #3. Please see Figure 7 for sample collected images for each sequence.

It's worthy noting that segmentation is crucial for marker detection and pose estimation for all marker systems. Thus, for the fair comparison, we fine tune the segmentation parameters for each marker algorithm unless it already uses advanced approaches like adaptive thresholding, line detection, etc. Specifically, we use threshold of 60 instead of default 100 for ARToolKit, 15 and 2 for `AdaptiveThresholdWindowSize` and `AdaptiveThresWindowSize_range` instead of default -1 and 0 for ArUco. Please refer to the supplementary material for the performance comparison between their default setups and our fine tuned versions.

## 5.2. Dictionary Size vs. Tracking Distance

Table 1 shows the comparison of dictionary size vs. tracking distance of different tag systems. Generally speaking, more tag bits offer more spaces to encode identities, but sacrifice tracking distance as region for each bit gets

Table 1: Dictionary size vs. tracking distance.

| Tag | Dictionary Size | Distance (m) |
|---|---|---|
| ARToolKit | 10[3] | 1.199 |
| ARToolKitPlus | 512 | 1.154 |
| ArUco (16h3) | 250 | 1.309 |
| ArUco (25h7) | 100 | 1.187 |
| ArUco (36h12) | 250 | 1.199 |
| RuneTag | 17,000 | 0.221 |
| ChromaTag | 30 | 0.560 |
| AprilTag (16h5) | 30 | 1.220 |
| AprilTag (25h7) | 242 | 1.171 |
| AprilTag (25h9) | 35 | 1.226 |
| AprilTag (36h9) | 5,329 | 1.223 |
| AprilTag (36h11) | 587 | 1.168 |
| TopoTag (3x3) | 128 | 1.204 |
| TopoTag (4x4) | 16,384 | 1.055 |
| TopoTag (5x5) | 8,388,608 | 0.670 |

Table 2: Detection accuracy (with run time).

| Tag | Recall (%) | Precision (%) | Time (ms) |
|---|---|---|---|
| ARToolKit | 99.990 | 99.880 | 5.864 |
| ARToolKitPlus | 98.297 | 100.000 | 9.314 |
| ArUco (16h3) | 100.000 | 99.910 | 54.319 |
| ArUco (25h7) | 99.009 | 100.000 | 53.930 |
| ArUco (36h12) | 99.470 | 100.000 | 56.001 |
| RuneTag | 0.281 | 100.000 | 455.832 |
| ChromaTag | 9.088 | 9.190 | 9.103 |
| AprilTag (16h5) | 77.285 | 99.883 | 246.762 |
| AprilTag (25h7) | 75.711 | 100.000 | 244.433 |
| AprilTag (25h9) | 80.405 | 100.000 | 251.275 |
| AprilTag (36h9) | 78.704 | 100.000 | 240.694 |
| AprilTag (36h11) | 77.577 | 100.000 | 241.314 |
| TopoTag | 100.000 | 100.000 | 33.638 |



Figure 8: Recall and precision by different points on Seq #3. Best viewed in color.

smaller. We can see that TopoTag achieves comparable tracking range when dictionary size is small (9-bit), while offers significant larger tracking range when dictionary size extends to tens of thousands (16-bit vs. RuneTag with the state-of-the-art most identities). In addition, TopoTag offers the scalability of extending dictionary size to millions with still acceptable tracking distance (25-bit).

### 5.3. Detection Accuracy

Table 2 summarizes the detection results and Figure 8 highlights the recall and precision of different captured points on Seq #3. We follow the metrics in [16]. True positives (TP) are defined as when the tag is correctly detected, i.e. locate the tag and correctly identify the ID (have at least 50% intersection over union between the detection and the ground truth). False positives (FP) are defined as detections that do not identify the location and ID correctly. False negatives (FN) are defined as any marker that is not identified correctly. Precision is $\frac{TP}{TP+FP}$ and recall is $\frac{TP}{TP+FN}$.

TopoTag performs perfectly on all three sequences, achieving 100% on both recall and precision. All marker systems except ChromaTag work great and achieve > 99.5% on precision due to their unique false positive rejection techniques. However, most systems except ARToolKit, ARToolKitPlus and ArUco fail to achieve a high recall, i.e. < 81%. Figure 8 shows that all previous systems degrade on recall or precision or both when markers are viewed from large angles. This probably results from large distortion, decreased lightness and blur issues, which distracts marker detection. TopoTag, on the other hand, has no obvious degradation. ChromaTag performs worse on both possibly due to the cluttered colored background and relative low brightness of collected images distracting its detection based on

color information. RuneTag performs the worst with lowest recall < 0.3% and fails to detect any frame on Seq #3.

**False Positive Rejection.** Since all our dataset contain valid tags. FP mainly focuses on the background excluding the tag regions. To better evaluate FP, as in [45], we further run the experiment on LabelMe [40] dataset, which consists of 207,883[4] images of natural scenes from a wide variety of indoor and outdoor environments, none of which contain any valid fiducial markers. We run this test for ARToolKit, ARToolKitPlus, ArUco, AprilTag and TopoTag as they achieve the state-of-the-art detect accuracy results on our dataset as shown in Table 2. There are 9756 false positives returned by ARToolKit, 49321 by AprilTag (16h5) and 146 by ArUco (16h3). As a contrast, TopoTag and ARToolKitPlus both have no false positives.

### 5.4. Localization Jitter and Accuracy

Localization jitter and accuracy are evaluated on Seq #3.

#### 5.4.1 Pose Error

We evaluate the accuracy between each point and its next point, where robot's measurements serve as the groundtruth.

---

[3] 10 tags are provided in ARToolKit package. Theoretically, any pattern can be used for tag design, but without providing a method.

[4] This is the latest LabelMe dataset size, which is slightly different the size of 180,829 from that was used in [32, 45].
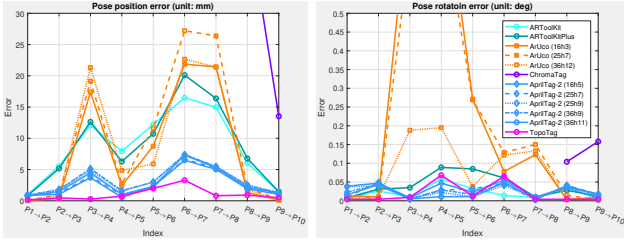
Figure 9: Pose position (left) and rotation (right) error comparison. We trim the figures for better visualization. Please refer to our supplementary material for full figures.
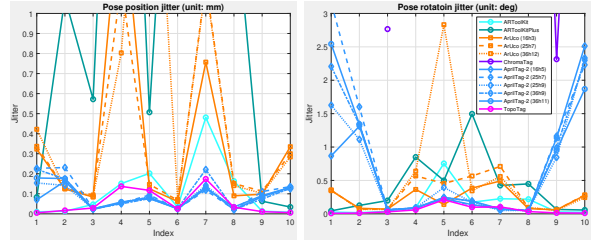


Figure 10: Pose position (left) and rotation (right) jitter comparison. We trim the figures for better visualization. Please refer to our supplementary material for full figures.

Table 3: Average and maximum pose errors for each tag. Best results are shown in bold.

| Tag | position (mm) | | rotation (deg) | |
|---|---|---|---|---|
| | avg | max | avg | max |
| ARToolKit | 8.639 | 16.499 | 0.022 | 0.058 |
| ARToolKitPlus | 8.923 | 20.101 | 0.040 | 0.089 |
| ArUco (16h3) | 8.191 | 21.876 | 0.248 | 0.908 |
| ArUco (25h7) | 10.049 | 27.212 | 0.225 | 0.765 |
| ArUco (36h12) | 8.768 | 22.663 | 0.078 | 0.195 |
| ChromaTag | 29.586 | 45.643 | 0.131 | 0.158 |
| AprilTag (16h5) | 2.894 | 7.287 | 0.031 | 0.055 |
| AprilTag (25h7) | 2.704 | 6.641 | 0.026 | 0.049 |
| AprilTag (25h9) | 3.178 | 7.320 | 0.024 | **0.041** |
| AprilTag (36h9) | 3.228 | 7.394 | 0.024 | 0.047 |
| AprilTag (36h11) | 2.549 | 6.489 | 0.022 | 0.049 |
| TopoTag | **1.011** | **3.289** | **0.019** | 0.068 |

Table 4: Average and maximum pose jitters for each tag. Best results are shown in bold.

| Tag | position (mm) | | rotation (deg) | |
|---|---|---|---|---|
| | avg | max | avg | max |
| ARToolKit | 0.112 | 0.481 | 0.160 | 0.754 |
| ARToolKitPlus | 1.134 | 3.584 | 0.421 | 1.496 |
| ArUco (16h3) | 0.363 | 1.636 | 0.230 | 0.491 |
| ArUco (25h7) | 0.364 | 1.155 | 0.322 | 0.710 |
| ArUco (36h12) | 0.573 | 2.553 | 0.526 | 2.832 |
| ChromaTag | 49.880 | 130.958 | 8.479 | 14.616 |
| AprilTag (16h5) | 0.079 | 0.163 | 0.654 | 2.512 |
| AprilTag (25h7) | 0.104 | 0.231 | 0.879 | 3.160 |
| AprilTag (25h9) | 0.087 | **0.154** | 0.673 | 2.333 |
| AprilTag (36h9) | 0.102 | 0.222 | 0.753 | 2.299 |
| AprilTag (36h11) | 0.092 | 0.179 | 0.734 | 2.543 |
| TopoTag | **0.055** | 0.173 | **0.058** | **0.211** |

Since there are in total 10 points in Seq #3, 9 accuracy values will be computed. See Figure 9 for the results of both position and rotation accuracies. Average and maximum pose errors for each tag can be seen in Table 3. TopoTag outperforms all previous systems in position error by a large margin (over 60% error reduction by average and about 50% by max compared to the 2nd best) and is comparable with the state-of-the-art on rotation error.

### 5.4.2 Pose Jitter

Jitters are evaluated at each point using the standard deviation (STD) metric. See Figure 10 for the result. Average and maximum jitter for each tag can be seen in Table 4. TopoTag outperforms all previous systems in rotation jitter by a large margin (over 63% jitter reduction by average and over 57% by max compared to the 2nd best) and is comparable with the state-of-the-art on position jitter.

### 5.4.3 Vertex Jitter

Vertex jitter measures the noise of the 2D feature point estimation, whose errors will propagate to the estimation of 6-

DoF pose. We compare with two best previous methods, i.e. AprilTag and ArUco. Both AprilTag and ArUco are square markers using intersections of quad lines to achieve sub-pixel vertex precision. RUNE-Tag and ChromaTag are not evaluated as they fail to reliably detect all positions in Seq #3, i.e. number of detected frames for a point is less than 50[5]. Square markers, like ARToolKitPlus and ChromaTag, theoretically will have similar performance as AprilTag and ArUco. ARToolKit is not evaluated as it uses correlation against a database to detect instead of finding fixed corners. All candidate methods are evaluated with 16-bit tags (i.e. AprilTag's 16h5 and ArUco's 16h3). Similar to pose jitter evaluation, STD metric is used.

Results can be seen in Figure 11. TopoTag performs consistently the best or comparable to the state-of-the-art across all points, especially when the marker angles become larger (e.g. $\geq 60°$) and with image blur (see P1, P2, P9 and P10). AprilTag performs better than ArUco when marker angles are relative small ($\leq 30°$, see P3-P8) thanks to its edge refinement, but become worse for larger marker angles.

---

[5]ChromaTag fails to reliably detect P2, P4, P5, P6 and P7; and all 10 positions are failed for RUNE-Tag.
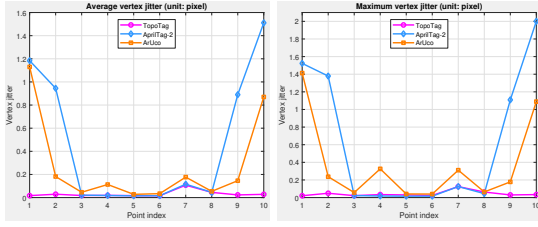
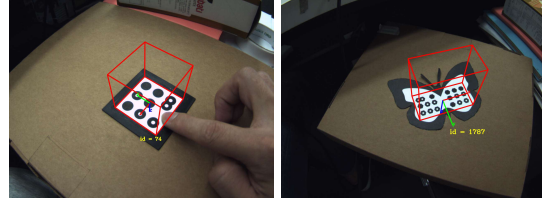Figure 11: Average and maximum vertex jitter comparison.



Figure 12: Detection and pose estimation of two customized TopoTags. Best viewed in color.
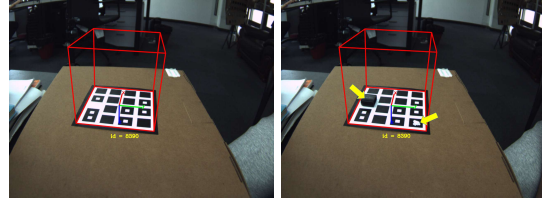


Figure 13: TopoTag detection under occlusion indicated by yellow arrows (before and after). Best viewed in color.
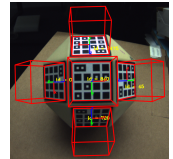
Table 5: Pose estimation for different bits and shapes.

| Different Bits & Shapes | Pose Accuracy | | | | Pose Jitter | | | |
|---|---|---|---|---|---|---|---|---|
| | position (mm) | | rotation (deg) | | position (mm) | | rotation (deg) | |
| | avg | max | avg | max | avg | max | avg | max |
| 3x3, circle | 1.073 | 3.299 | 0.016 | 0.048 | 0.069 | 0.265 | 0.080 | 0.205 |
| 3x3, square | 0.837 | 2.780 | 0.022 | 0.065 | 0.085 | 0.383 | 0.081 | 0.276 |
| 4x4, circle | 0.995 | 2.867 | 0.017 | 0.057 | 0.058 | 0.192 | 0.073 | 0.272 |
| 4x4, square | 1.011 | 3.289 | 0.019 | 0.068 | 0.055 | 0.173 | 0.058 | 0.211 |

## 5.5. Speed

### 5.5.1 Dictionary Computation

Dictionary computation usually is a time-consuming operation due to the specially designed lexicode generation algorithm and Hamming distance strategy for high robustness. Though there is no need to do this online, it's still meaningful to make it efficient enough. ArUco takes about 8, 20 and 90 minutes respectively for dictionaries of sizes 10, 100 and 1000 [20], while it can take several days to generate 36-bit tags for AprilTag [32]. As TopoTag supports full tag bits for encoding, it takes almost no time as ID can be directly mapped to the binary code string for dictionary generation.

### 5.5.2 Tag Detection

Table 2 (last column) shows the running time comparison. TopoTag takes less time than ArUco (38% ⇓), AprilTag (86% ⇓) and RuneTag (93% ⇓). Though ARToolKit, ChromaTag and ARToolKitPlus run faster, they offer significantly less identities, see Table 1. For TopoTag, most time is spent on segmentation (68.8%), followed by decoding & vertex estimation (29.7%) and pose estimation (1.5%).

No parallelization is utilized in current TopoTag implementation, which will normally bring further speed gain otherwise. We have implemented the segmentation part in a single pipeline on a Lattice FPGA (LFE5UM-45 with 44k LUTs, 1.9 Mb RAM and without using external DDR), which is decreased to $< 100$ us achieving $230\times$ speedup.

## 5.6. Flexible Shape Support

TopoTag supports both customized outer and inner shapes as long as topological structure is maintained. Figure 1 shows three TopoTags with various inter shapes like squares, circles, hexagons and different outer shapes like squares and customized like a bufferfly. Figure 12 shows our algorithm running for these customized TopoTags.

Experiments show that tags with different shapes have compatible results. Please see Table 5 for detailed comparison. It's worthy noting that all these four different TopoTags have 100% on both detection recall and precision which further validates the robustness of our system.

## 5.7. Occlusion Support

TopoTag can handle occlusion as long as topological structure is preserved. See the left image of Figure 12 for an example working under occlusion. We can further incorporate Hamming distance in tag design to handle occlusions where topological is no longer maintained by sacrificing dictionary size. Figure 13 shows an example of working under occlusion by incorporating Hamming distance of 3, where dictionary size will decrease from 16384 to 1024.

To handle more severe occlusions, similar to [1, 20], we can use multiple tags in a grid to increase the probability of seeing complete markers. Other forms can be also considered. On right, we show an example of achieving 360°-freedom tracking via using 18 TopoTags on a rhombicuboctahedron-shaped object.



## 6. Conclusions

We present TopoTag, a new topological-based fiducial marker and detection algorithm that utilizes topological information to achieve high robustness, achieving perfect de-

tection accuracy on the large dataset we collected and also datasets from others. We show that all tag bits can encode identities without sacrificing detection accuracy, thus achieving rich identification and scalability. TopoTag offers more feature correspondences for a better pose estimation. TopoTag achieves the best performance in localization jitter and accuracy, and at the same time supports occlusion and flexible shapes. We also collect a large dataset of TopoTag and six previous state-of-the-art tags for future evaluation, involving in-plane and out-of-plane rotations, image blur, various distances and cluttered background, etc.

## References

[1] B. Atcheson, F. Heide, and W. Heidrich. CALTag: High Precision Fiducial Markers for Camera Calibration. *Int. Workshop on Vision, Modeling and Visualization (VMV)*, 2010. 1, 2, 8

[2] R. Bencina and M. Kaltenbrunner. The design and evolution of fiducials for the reactivision system. In *Proceedings of the Third International Conference on Generative Systems in the Electronic Arts*, 2005. 2

[3] R. Bencina, M. Kaltenbrunner, and S. Jorda. Improved topological fiducial tracking in the reactivision system. In *2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'05)-Workshops*, pages 99–99. IEEE, 2005. 2

[4] F. Bergamasco, A. Albarelli, L. Cosmo, E. Rodola, and A. Torsello. An accurate and robust artificial marker based on cyclic codes. *IEEE transactions on pattern analysis and machine intelligence*, 38(12):2359–2373, 2016. 2

[5] F. Bergamasco, A. Albarelli, E. Rodolà, and A. Torsello. RUNE-Tag: A high accuracy fiducial marker with strong occlusion resilience. *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pages 113–120, 2011. 1, 2, 4, 5

[6] F. Bergamasco, A. Albarelli, and A. Torsello. Pi-tag: a fast image-space marker design based on projective invariants. *Machine vision and applications*, 24(6):1295–1310, 2013. 2

[7] L. Calvet, P. Gurdjos, and V. Charvillat. Camera tracking using concentric circle markers: Paradigms and algorithms. In *2012 19th IEEE International Conference on Image Processing*, pages 1361–1364. IEEE, 2012. 2

[8] L. Calvet, P. Gurdjos, C. Griwodz, and S. Gasparini. Detection and accurate localization of circular fiducials under highly challenging conditions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 562–570, 2016. 1, 2, 4

[9] Y. Cho, J. Lee, and U. Neumann. A multi-ring color fiducial system and an intensity-invariant detection method for scalable fiducial-tracking augmented reality. In *IWAR*, 1998. 1, 2

[10] D. Claus and A. W. Fitzgibbon. Reliable fiducial detection in natural scenes. In *European Conference on Computer Vision*, pages 469–480. Springer, 2004. 2

[11] D. Claus and A. W. Fitzgibbon. Reliable automatic calibration of a marker-based position tracking system. In *2005*

[12] E. Costanza. D-touch: A consumer-grade tangible interface module and musical applications. In *Proceedings of Conference on HumanComputer Interaction*, 2003. 2

[13] E. Costanza and J. Robinson. A Region Adjacency Tree Approach to the Detection and Design of Fiducials. In *Video Vision and Graphics*, pages 63–69, 2003. 2

[14] H. Cruz-Hernández and L. G. de la Fraga. A fiducial tag invariant to rotation, translation, and perspective transformations. *Pattern Recognition*, 81:213–223, 2018. 1

[15] V. F. da Camara Neto, D. B. de Mesquita, R. F. Garcia, and M. F. M. Campos. On the design and evaluation of a precise scalable fiducial marker framework. In *2010 23rd SIBGRAPI Conference on Graphics, Patterns and Images*, pages 216–223. IEEE, 2010. 1, 2

[16] J. DeGol, T. Bretl, and D. Hoiem. Chromatag: a colored marker and fast detection algorithm. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1472–1481, 2017. 1, 2, 3, 4, 5, 6

[17] M. Fiala. Artag, a fiducial marker system using digital techniques. In *2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'05)*, volume 2, pages 590–596. IEEE, 2005. 1, 2, 3

[18] M. Fiala. Designing highly reliable fiducial markers. *IEEE Transactions on Pattern analysis and machine intelligence*, 32(7):1317–1324, 2010. 1

[19] D. Flohr and J. Fischer. A lightweight id-based extension for marker tracking systems. In *Eurographics Symposium on Virtual Environments (EGVE) Short Paper Proceedings*, pages 59–64, 2007. 1, 2

[20] S. Garrido-Jurado, R. Muñoz-Salinas, F. J. Madrid-Cuevas, and M. J. Marín-Jiménez. Automatic generation and detection of highly reliable fiducial markers under occlusion. *Pattern Recognition*, 47(6):2280–2292, 2014. 1, 2, 3, 5, 8

[21] S. Garrido-Jurado, R. Munoz-Salinas, F. J. Madrid-Cuevas, and R. Medina-Carnicer. Generation of fiducial marker dictionaries using mixed integer linear programming. *Pattern Recognition*, 51:481–491, 2016. 1, 2, 3

[22] L. B. Gatrell, W. A. Hoff, and C. W. Sklair. Robust image features: Concentric contrasting circles and their image extraction. In *Cooperative Intelligent Robotics in Space II*, volume 1612, pages 235–245. International Society for Optics and Photonics, 1992. 2

[23] R. Hartley and A. Zisserman. *Multiple view geometry in computer vision*. Cambridge university press, 2003. 4

[24] M. Kaltenbrunner and R. Bencina. reactivision: a computer-vision framework for table-based tangible interaction. In *Proceedings of the 1st international conference on Tangible and embedded interaction*, pages 69–74. ACM, 2007. 2

[25] H. Kato and M. Billinghurst. Marker tracking and hmd calibration for a video-based augmented reality conferencing system. In *Proceedings 2nd IEEE and ACM International Workshop on Augmented Reality (IWAR'99)*, pages 85–94. IEEE, 1999. 2, 5

Seventh IEEE Workshops on Applications of Computer Vision (WACV/MOTION'05)-Volume 1, volume 1, pages 300–305. IEEE, 2005. 2

[26] C. N. Klokmose, J. B. Kristensen, R. Bagge, and K. Halskov. Bullseye: high-precision fiducial tracking for table-based tangible interaction. In *Proceedings of the Ninth ACM International Conference on Interactive Tabletops and Surfaces*, pages 269–278. ACM, 2014. 2

[27] V. A. Knyaz. The development of new coded targets for automated point identification and non-contact 3d surface measurements. *IAPRS*, 5:80–85, 1998. 2

[28] V. Lepetit and P. Fua. Keypoint recognition using randomized trees. *IEEE transactions on pattern analysis and machine intelligence*, 28(9):1465–1479, 2006. 2

[29] D. W. Marquardt. An algorithm for least-squares estimation of nonlinear parameters. *Journal of the society for Industrial and Applied Mathematics*, 11(2):431–441, 1963. 4

[30] J. Molineros and R. Sharma. Real-time tracking of multiple objects using fiducials for augmented reality. *Real-Time Imaging*, 7(6):495–506, 2001. 1

[31] L. Naimark and E. Foxlin. Circular data matrix fiducial system and robust image processing for a wearable vision-inertial self-tracker. In *Proceedings of the 1st International Symposium on Mixed and Augmented Reality*, page 27. IEEE Computer Society, 2002. 2

[32] E. Olson. Apriltag: A robust and flexible visual fiducial system. In *2011 IEEE International Conference on Robotics and Automation*, pages 3400–3407. IEEE, 2011. 1, 2, 3, 4, 6, 8

[33] C. B. Owen, F. Xiao, and P. Middlin. What is the best fiducial? In *The First IEEE International Workshop Agumented Reality Toolkit,*, pages 8–pp. IEEE, 2002. 2, 4

[34] M. Ozuysal, M. Calonder, V. Lepetit, and P. Fua. Fast keypoint recognition using random ferns. *IEEE transactions on pattern analysis and machine intelligence*, 32(3):448–461, 2010. 2

[35] M. G. Prasad, S. Chandran, and M. S. Brown. A motion blur resilient fiducial for quadcopter imaging. In *2015 IEEE Winter Conference on Applications of Computer Vision*, pages 254–261. IEEE, 2015. 1, 2

[36] J. Rekimoto. Matrix: A realtime object identification and registration method for augmented reality. In *Proceedings. 3rd Asia Pacific Computer Human Interaction (Cat. No. 98EX110)*, pages 63–68. IEEE, 1998. 2

[37] J. Rekimoto and Y. Ayatsuka. Cybercode: designing augmented reality environments with visual tags. In *Proceedings of DARE 2000 on Designing augmented reality environments*, pages 1–10. ACM, 2000. 2

[38] M. Rohs and B. Gfeller. Using camera-equipped mobile phones for interacting with real-world objects. *Advances in pervasive computing*, 176:265–271, 2004. 2

[39] F. J. Romero-Ramirez, R. Muñoz-Salinas, and R. Medina-Carnicer. Speeded up detection of squared fiducial markers. *Image and Vision Computing*, 76:38–47, 2018. 1

[40] B. C. Russell, A. Torralba, K. P. Murphy, and W. T. Freeman. Labelme: a database and web-based tool for image annotation. *International journal of computer vision*, 77(1-3):157–173, 2008. 4, 6

[41] J. Sattar, E. Bourque, P. Giguere, and G. Dudek. Fourier tags: Smoothly degradable fiducial markers for use in human-robot interaction. In *Fourth Canadian Conference on Computer and Robot Vision (CRV'07)*, pages 165–174. IEEE, 2007. 2

[42] F. Schweiger, B. Zeisl, P. Georgel, G. Schroth, E. Steinbach, and N. Navab. Maximum detector response markers for sift and surf. In *Vision, Modeling and Visualization Workshop (VMV)*, 2009. 2

[43] K. Tateno, I. Kitahara, and Y. Ohta. A nested marker for augmented reality. In *2007 IEEE Virtual Reality Conference*, pages 259–262. IEEE, 2007. 2

[44] D. WAGNER. Artoolkitplus for pose tracking on mobile devices. In *Proceedings of 12th Computer Vision Winter Workshop (CVWW'07), February*, 2007. 2, 3, 4, 5

[45] J. Wang and E. Olson. Apriltag 2: Efficient and robust fiducial detection. In *2016 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 4193–4198. IEEE, 2016. 1, 2, 3, 4, 5, 6

[46] A. Xu and G. Dudek. Fourier tag: A smoothly degradable fiducial marker system with configurable payload capacity. In *2011 Canadian Conference on Computer and Robot Vision*, pages 40–47. IEEE, 2011. 2