

AprilTrack: A Particle-Filter based AprilTag Tracker

Daniel Pfrommer

Abstract

The AprilTrack algorithm proposed in this paper seeks to bridge two existing fields of research in computer vision literature: fiducial detection and monocular object tracking. With a rise in the popularity of fiducials, or artificial markers, as a source of ground-truth pose estimation for visual odometry datasets and mobile robots, there has been considerable interest in designing robust and blur-resistant fiducials. Rather than devising a new fiducial, this work attempts to improve on existing fiducials, such as the popular AprilTag marker, by augmenting fiducial detection algorithms with a particle-filter based monocular object tracker which reliably tracks a fiducial's position and orientation in six degrees of freedom, allowing for tag-based localization in situations where current detectors would be unable to do so. This novel patch-based particle filter tracking algorithm, which is extended to simultaneously track multiple tags in ego-motion scenarios, is qualitatively and quantitatively evaluated on several AprilTag sequences, demonstrating that it is capable of reliably localizing blurred tags.

1. Introduction

Visual tracking of fiducial markers (or simply fiducials) plays an important role in many robotics and computer vision related fields and applications, including augmented and virtual reality [2], visual odometry/SLAM [3], and object tracking. However, as many of these applications involve the use of fiducials in structured indoors environments where illumination is often limited, long exposure times to compensate for the reduced lighting produce blurred images, especially when taken from a mobile platform such as a quadrotor or a handheld device.

The focus of this work is primarily on providing better pose estimation for visual odometry datasets which use AprilTags placed to known locations to provide the ground truth pose data. As the pose estimation in this scenario can be done in a post-processing step after the data has been collected, the current implementation of AprilTrack is not designed to run in realtime.

We intend to explore the use of AprilTrack in realtime scenarios in future research.

1.1. Previous Work

The majority of previous literature in fiducial detection has been focused single-frame fiducial detection [1, 7, 5]. A common assumption in these algorithms is that the image being processed shows little motion blur, an expectation which breaks down in empirical situations. For instance, the widely-used AprilTag marker [5] and its successor, the AprilTag 2 [8], rely on strong gradients between the tag's background and the surrounding white border to detect the tag. As shown in figure 1, oftentimes such algorithms fail to detect tags that have been blurred, either missing them entirely or mistakenly discarding them when eliminating false positives through some encoded fiducial ID (for AprilTags this consists of a 6x6 binary grid in the center of the tag which is used to verify the validity of the detection).

Previously, there has been interest in designing blur resilient fiducials. Meghshyam et al. [7] have developed ringed markers to ensure that a region of the tag will always be perpendicular to the blur direction. However, this approach, while resistant to linear blur, does not take into account more complex blurs (e.g resulting from the tag pivoting around the vertical axis). In addition, the increased resilience to blur comes at the cost of localization accuracy as only the center point of the ringed tag can be reliably determined [7], whereas with other fiducials, such as the AprilTag [5], full pose estimation in 6 degrees of freedom (DOF) can be done based on just a single tag. Another significant disadvantage of the ringed marker from Meghshyam et al. is the limited number of tag ID's, as their proposed encoding method allows for a maximum of 4 distinct tags [7].

Unlike previous attempts to engineer more resilient tags, the algorithm proposed here exploits the temporal coherence between successive frames in a video sequence to track a tag's location between successful detections. As many high-blur scenarios involve a continuous stream of images taken at equidistant time intervals (e.g a quadrotor-mounted moving camera viewing tags or a stationary camera tracking an AprilTag-

labeled moving object), we propose to use knowledge of a fiducial’s prior location to augment the detection algorithm. As such, this paper presents a novel particle-filter-based tracking algorithm (AprilTrack) capable of tracking AprilTags in 6 DOF with both joint-tag motion and camera ego-motion.

Previous work involving the tracking of blurred objects includes the Blur-driven Tracker (BLUT) framework proposed by Wu et al. [9] for tracking motion-blurred targets. BLUT, a particle-filter based tracker, works on the observation that blur strength and direction can be used estimate an object’s velocity and provide valuable information for successfully tracking a blurred object. While capable of reliably tracking blurred objects, BLUT suffers from the disadvantage that it can do so only in two dimensions (i.e it is assumed that the object always faces the camera) and so cannot track a sequence of involving arbitrary tag rotations and translations.

The novel method proposed in this paper augments the AprilTag detection algorithm originally designed by Olson [5] and improved by Wang et al [8] with a particle-filter based algorithm to reliably track an AprilTag through an arbitrary sequence of images, providing valuable fiducial-based position information when existing methods would be unable to do so. In addition, the proposed algorithm is extended to incorporate the joint motion of an arbitrary number of tags to increase the accuracy of the tracked orientation and position.

Overview. The remainder of the paper is organized as follows. Section 2 gives an overview of the method proposed in this paper and describes the algorithm in detail. Section 3 describes the implementation of the algorithm and evaluates its performance on a manually-labeled dataset of blurred moving tags both quantitatively and qualitatively. In section 4 we discuss the lim-



Figure 1. A blurred AprilTag. Such tags are difficult for still frame based algorithms to detect as the payload (the ID encoded in the center of the tag) will cause the tag to be discarded as a false positive.

itations of this approach and suggest possible avenues of future research.

2. Method

The AprilTrack algorithm is similar in concept to the BLUT tracker by Wu et al. [9] in that it uses a particle filter to integrate all available information about a marker’s position. Overall, the description of the algorithm can be broken down into three parts: one describing the particle filter, another on the method for evaluating candidate particles, and one describing the modifications made to simultaneously track multiple tags.

2.1. The Particle Filter

The particle filter used for the AprilTrack algorithm is a Bayesian Sequential Importance Resampling (SIR) technique which approximates a posterior distribution of state variables at discrete time intervals. There are three major steps in the algorithm: prediction, update, and resampling. We use the vector x_t to denote the 13 dimensional state vector describing the position, orientation, velocity, and rotational velocity of the tracked target at time t

$$x_t = \begin{bmatrix} r \\ \dot{r} \\ q \\ \omega \end{bmatrix} \quad (1)$$

where r is the three dimensional position of the tracked object, \dot{r} is the three dimensional velocity, q is a four dimensional unit quaternion containing the pose of the tracked object, and ω is an angular velocity with components ω_x , ω_y , and ω_z .

Optimal filtering. The prediction step can be formulated in a Bayesian setting as the computation of the prior distribution $p(x_t|z_{t-1})$ from the filtering distribution $p(x_{t-1}|z_{1:t-1})$ where z_t represents the measurement (appearance of a fiducial) taken at time t . As in [6]:

$$p(x_t|z_{1:t-1}) = \int p(x_t|x_{t-1})p(x_{t-1}|z_{1:t-1})dx_{t-1} \quad (2)$$

Likewise, the update step can be formulated in a similar manner, where the prior distribution calculated in the prediction step is updated with the latest measurement to obtain the posterior over x_t . Mathematically this can be derived using Bayes’ law to express $p(x_t|z_{1:t})$ in terms of $p(z_t|x_t)$ and $p(x_t|z_{1:t-1})$. Using Bayes’ law yields

$$p(x_t|z_{1:t}) = p(x_t|z_t, z_{1:t-1})$$

$$p(x_t|z_t, z_{1:t-1}) = \frac{p(z_t|x_t, z_{1:t-1})p(x_t|z_{1:t-1})p(z_{1:t-1})}{p(z_{1:t})}$$

As the model $p(x_t|x_{t-1})$ is Markovian, the current measurement is conditionally independent of previous measurements given x_t , meaning that $p(z_t|x_t, z_{1:t-1})$ can be written as $p(z_t|x_t)$. $p(x_t|z_{1:t})$ can also be expressed as $p(z_t|z_{1:t-1})p(x_t|z_{1:t-1})$ resulting in the final update equation [6]

$$p(x_t|z_{1:t}) = \frac{p(z_t|x_t)p(x_t|z_{1:t-1})}{p(z_t|z_{1:t-1})} \quad (3)$$

Particle Filters. In this regard, particle filters attempt to approximate the posterior distribution at time $t-1$, $p(x_{0:t-1}|z_{1:t-1})$, using Monte-Carlo methods. A particle filter consists of a set of N samples $\{(x_{t-1}^i, w_{t-1}^i)\}_{i=1}^N$ drawn from the proposal distribution $q(x^i)$ where each sample is called a *particle*. This allows for the filtering distribution $p(x_{0:t-1}|z_{1:t-1})$ to be approximated by (as described in [6])

$$p(x_{0:t-1}|z_{1:t-1}) \approx \sum_{i=1}^N w_{t-1}^i \delta_{x_{0:t-1}^i} \quad (4)$$

where

$$\sum_{i=1}^N w_t^i = 1 \quad (5)$$

As the particles are drawn from the proposal distribution $q(x)$, the sample weight w^i must take into account any discrepancy between $p(x)$ and $q(x)$ so that $w_t^i \propto p(x_{1:t}^i|z_{1:t})/q(x_{1:t}^i|z_{1:t})$. From this observation we note that therefore:

$$w_t^i \propto \frac{p(x_{0:t}^i|z_{1:t})}{q(x_{0:t}^i|z_{1:t})}$$

In the simplified case of a Sequential Importance Resampling algorithm, the prediction and update equations derived from (2) and (3) can, in the context of a SIR particle filter, be simplified to:

$$x_t^i \sim p(x_t|x_{t-1}^i) \quad (6)$$

$$w \propto p(z_t|x_t^i) \quad (7)$$

where $x_t^i \sim p(x_t|x_{t-1}^i)$ denotes that, during the prediction step, the particles for the prior distribution at time t are sampled from the known distribution $p(x_t|x_{t-1}^i)$. The distribution $p(z_t|x_t^i)$ for updating particle weights will be defined in section 2.2.

In AprilTrack, we model $p(x_t|x_{t-1}^i)$ as a multivariate gaussian distribution with standard deviation σ

around a mean μ where

$$\mu = \begin{bmatrix} r_{t-1}^i + \Delta t \dot{r}_{t-1}^i \\ \dot{r}_{t-1}^i \\ q_{t-1}^i \dot{q}_{t-1}^i \\ \omega_{t-1}^i \end{bmatrix} \quad (8)$$

In equation (8), Δt is the elapsed time and \dot{q}_{t-1}^i is the quaternion representation of $\Delta t \omega_{k-1}^i$, as described in [4]. We also use the *noise quaternion* formulated in [4] to apply the gaussian noise to the quaternion $q_{k-1}^i \dot{q}_{k-1}^i$. For simplicity, we make the assumption that the input images have been taken at equidistant time intervals and incorporate Δt into the ω and \dot{r} components of σ .

After completing the update prediction and update steps and selecting the particle with the highest weight as the most likely candidate, the final resampling step simply consists of drawing a new set of N particles with equal weights (normalized to a total of one) from the distribution $p(x_t|z_{1:t})$ where $p(x_t|z_{1:t})$ is approximated by the set of weighted particles.

$$p(x_t|z_{1:t}) \approx \sum_{i=1}^N w_t^i \delta_{x_t^i} \quad (9)$$

2.2. Evaluation of Candidate Particles

The tracking of AprilTags in 6 DOF requires a method for evaluating the likelihood of candidate particles at arbitrary three dimensional positions and rotations, or $p(z_t|x_t^i)$. To this effect, the proposed method uses homographies to generate image *patches*, which are unprojected representations of an AprilTag for a particular pose and image.

Given the function $I([\lambda x, \lambda y, \lambda]^T)$ representing the intensity of pixel at a particular (x, y) of an image, a patch in the image can be defined by the function

$$P(x_p, y_p|H, I) = I(H \begin{bmatrix} x_p & y_p & 1 \end{bmatrix}^T) \quad (10)$$

where x_p and y_p are between -1 and 1 and H is a 3x3 matrix (a homography) mapping from $(x_p, y_p, 1)$ to $(\lambda x, \lambda y, \lambda)$ which incorporates the pose of an AprilTag relative to the camera.

This allows for the evaluation of candidate particle x^i by the measuring the similarity between the particle's patch $P^i(x_p, y_p|H^i, I)$ in the current image and the patch from the AprilTag's last successful detection, $P_{ref}(x_p, y_p|H_{ref}, I_{ref})$. Here, H^i can be calculated with

$$H^i = \begin{bmatrix} \hat{q}^i & r^i \end{bmatrix} \begin{bmatrix} \tau/2 & 0 & 0 \\ 0 & \tau/2 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad (11)$$

where \hat{q}^i is the rotation matrix representation of particle x^i 's rotation quaternion q , r is the particle's translation, and τ is the tag's edge length. The patch for the detected AprilTag uses the homography H_{ref} supplied by the AprilTag library, which computes the homography using the Direct Linear Transform algorithm [5].

To compare the patches efficiently, P^i and P_{ref} are evaluated at a regular grid from $(-1, -1)$ to $(1, 1)$ at intervals of $2/\rho$. This results in two $\rho \times \rho$ matrices M^i and M_{ref} . To compare the patches, the error between the two matrices is calculated using the correlation coefficient between the individual elements of the matrices.

$$\varepsilon(A, B) = \frac{1}{2} - \frac{\sum_{i=1}^{\rho} \sum_{j=1}^{\rho} (A_{ij} - \mu_A)(B_{ij} - \mu_B)}{2\sigma_A\sigma_B} \quad (12)$$

where A and B are the matrix approximations of two patches, σ_A and σ_B are the standard deviations of the respective matrices' elements and μ_A and μ_B are their means. By normalizing by mean and standard deviation, using correlation as a similarity metric (as opposed to the sum of the absolute difference or squared difference) ensures lighting invariance for comparing patches containing the same marker under different illuminations.

As in BLUT [9], the patch error $\varepsilon(M^i, M_{ref})$ is used to compute the observation likelihood $p(z_t | x_t^i)$ to update the weights and select the best particle with

$$p(z_t | x_t^i) = \exp(-\gamma \varepsilon(M^i, M_{ref})) \quad (13)$$

where the parameter γ allows for more or less selective resampling.

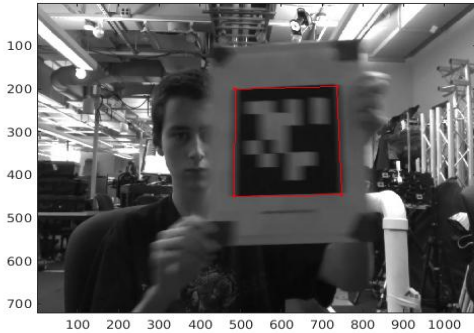


Figure 2. A blurred tag being tracked. The red border indicates that the tag was not detected with the AprilTag detector.

2.3. Multi-Tag Scenarios

To adapt this method to multi-tag scenarios, several extra parameters have to be introduced. In order to jointly track multiple tags, it is assumed that the relative poses of the tags are known. For each tag T_k with ID k , the location x_k and orientation quaternion g_k are expressed such that

$$X_c = \hat{q}^i(\hat{g}_k X_t + x_k) + r^i \quad (14)$$

where X_t is a coordinate relative to the tag's reference frame, X_c is relative to the camera reference frame, and q^i and r^i are the position and orientation of particle x_t^i . In the simple case of jointly moving tags, the H^i from equation (11) can then be rewritten as

$$H_k^i = \begin{bmatrix} \hat{q}^i \hat{g}_k & \hat{q}^i x_k + r^i \end{bmatrix} \begin{bmatrix} \tau_k/2 & 0 & 0 \\ 0 & \tau_k/2 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad (15)$$

where H_k^i represents the homography for particle i and tag k and τ_k is the edge length of tag k .

In the ego-motion scenario, each particle's rotation and translation can be taken to be the camera's rotation and translation such that

$$\hat{q}^i X_c + r^i = \hat{g}_k X_t + x_k \quad (16)$$

resulting in the homography H_k^i and the corresponding patch approximation M_k^i where

$$H_k^i = \begin{bmatrix} \hat{q}^{i-1} \hat{g}_k & \hat{q}^{i-1} (x_k - r^i) \end{bmatrix} \begin{bmatrix} \tau_k/2 & 0 & 0 \\ 0 & \tau_k/2 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad (17)$$

In both ego-motion and joint tag motion, the observation likelihood $p(z_t | x_t)$ then becomes

$$p(z_t | x_t^i) = \exp\left(\sum_{k=1}^K -\gamma \varepsilon(M_k^i, M_{k,ref})\right) \quad (18)$$

where K is the number of tags being tracked and $M_{k,ref}$ is the reference patch approximation for tag k from the last successful detection.

3. Results and Discussion

The AprilTrack algorithm was primarily evaluated on a manually-labeled single-tag sequence with a moving tag. The algorithm was also run on an ego-

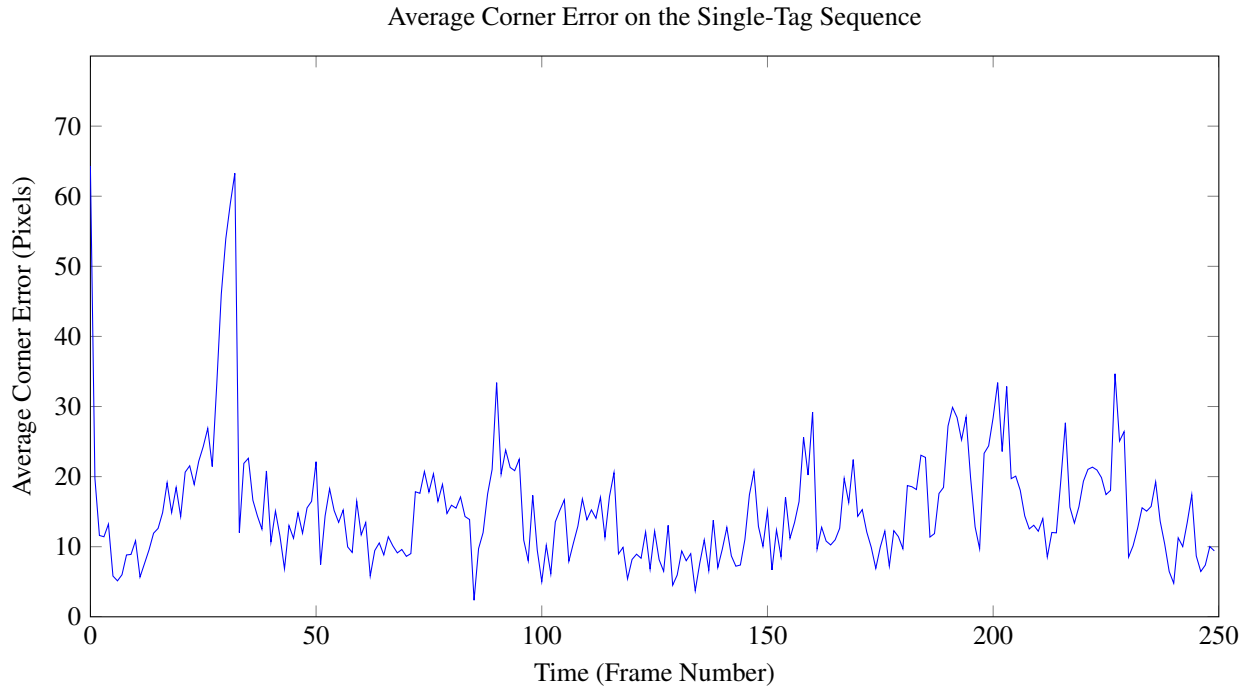


Figure 3. The average corner error between the hand-marked tag location and the tracked location. Much of the error comes from discrepancy between the tracker and the ground truth on the boundaries of the blurred tags. At around frames 30 and 200, the tracker loses the tag momentarily due to a sudden rotation but is able to recover without needing to be manually reset.

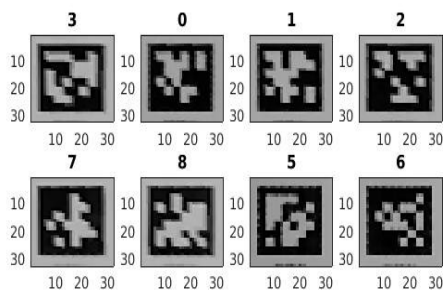


Figure 4. The reference patches of various tags (by ID) in a multi-Tag tracking scenario.

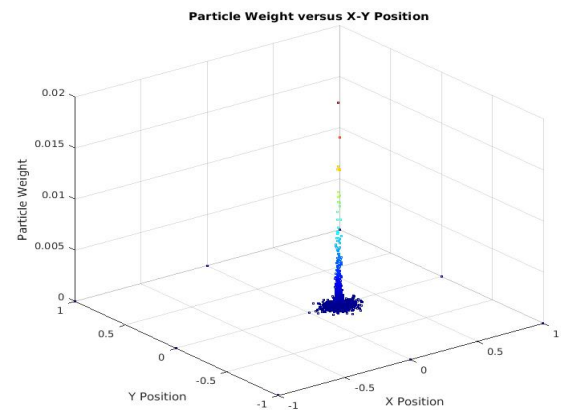


Figure 5. The distribution of particles and their weights with respect to the X,Y position of the particles.

motion multi-tag sequence which will be qualitatively discussed.¹

The specific implementation of the AprilTrack algorithm used in this experiment was written in MATLAB. For both scenarios, the parameters $\gamma = 10$, $\rho = 32$ were used. τ was 0.1936 m for all tags, slightly larger than the actual tag size of 0.1635 m to include some of the whitespace around the tag. In the single-tag scenario, the standard deviation σ of the gaussian noise used in the state transition distribution $p(x_t|x_{t-1}^i)$ was 0.01 for the r component, 0.05 for q , 0.02 for \dot{r} , and 0 for ω .

When evaluated on the single-tag sequence, the algorithm ran with $N = 3000$ at around 1 fps on an i7 Haswell Intel CPU. In experiments with fewer particles, it was found that the algorithm could track a tag well with as few as 1000 particles with around 5 fps. With a more optimized C++ implementation, the framerate could perhaps be brought up to 10 – 15 fps, making the algorithm (run with less aggressive parameters) a possible candidate for real-time scenarios. Currently, the MATLAB implementation is intended for offline post-processing of a sequence and so real-time performance was not a focus of this research.

As shown in figure 2.2, the pixel corner error hovers at around 10 pixels for most of the sequence. The tracker loses the tag momentarily around frames 30 and 200 but is able to recover rapidly without any manual intervention. Much of the error observed throughout the sequence is the result of disagreement on the boundary of a blurred tag rather than a systematic issue with the tracker itself.

When running with a multi-tag configuration, the tracker is able to make use of the additional information to more accurately track the camera's location than in a single-tag scenario. The main challenge for running the algorithm with multiple tags was not successfully tracking the tags—the algorithm was able to do that with half the particles of the single-tag scenario. Rather, the majority of the reprojection error arises from error in the relative tag positions. This opens up new potential avenues of research, such as tag-based bundle adjustment, to alleviate these issues.

4. Conclusion

This paper has presented the AprilTrack algorithm augmenting the standard AprilTag detection algorithm with a novel tracker capable of tracking the tags in 6 DOF by using homographies to evaluate candidate par-

ticles. Experimental results show that the tracker is capable of accurately tracking both single tags and consolidating information from multiple tags in ego-motion and joint tag-motion scenarios. In future research we plan to investigate other possible error metrics for patch comparison and look to deploy AprilTrack in realtime scenarios.

References

- [1] F. Bergamasco, A. Albarelli, E. Rodola, and A. Torsello. Rune-tag: A high accuracy fiducial marker with strong occlusion resilience. In *Computer Vision and Pattern Recognition (CVPR), 2011 IEEE Conference on*, pages 113–120. IEEE, 2011.
- [2] 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 *In IWAR*. Citeseer, 1998.
- [3] S. Houben, D. Droschel, and S. Behnke. Joint 3D Laser and Visual Fiducial Marker based SLAM for a Micro Aerial Vehicle. 2016.
- [4] E. Kraft. A quaternion-based unscented kalman filter for orientation tracking. In *Proceedings of the Sixth International Conference of Information Fusion*, volume 1, pages 47–54, 2003.
- [5] E. Olson. AprilTag: A robust and flexible visual fiducial system. *Proceedings - IEEE International Conference on Robotics and Automation*, pages 3400–3407, 2011.
- [6] E. Orhan. Particle Filtering, 2012.
- [7] M. G. Prasad, S. Chandran, and M. S. Brown. A motion blur resilient fiducial for quadcopter imaging. *Proceedings - 2015 IEEE Winter Conference on Applications of Computer Vision, WACV 2015*, pages 254–261, 2015.
- [8] J. Wang and E. Olson. Apriltag 2: Efficient and robust fiducial detection. In *Intelligent Robots and Systems (IROS), 2016 IEEE/RSJ International Conference on*, pages 4193–4198. IEEE, 2016.
- [9] Y. Wu, H. Ling, J. Yu, F. Li, X. Mei, and E. Cheng. Blurred target tracking by blur-driven tracker. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1100–1107, 2011.

¹Videos of the sequences and the tracker's performance are available online at <https://www.youtube.com/watch?v=kV8z3st9rlk> and <https://www.youtube.com/watch?v=F5cFTLBIJ1s>