

# Homography-Based Adaptive Robot Visual Tracking With Camera Parameter Convergence

Beixian Lai<sup>2</sup>, Yongping Pan<sup>1</sup>, Zhiwen Li<sup>2</sup>, and Changyun Wen<sup>3</sup>

**Abstract**—Parametric uncertainties are ubiquitous in vision-based robotic systems. However, existing adaptive visual servoing methods can not simultaneously achieve six-degree-of-freedom (6-DoF) robot pose control in the three-dimensional (3D) space and accurate camera parameter estimation. This paper considers robot manipulators with eye-to-hand monocular cameras under unknown extrinsic parameters and proposes a passivity-based kinematic control method called homography-based visual servoing with composite learning (CL-HBVS) to achieve 6-DoF robot pose tracking. The convergence of camera extrinsic parameters is achieved by designing a composite learning mechanism without the stringent condition of persistent excitation, which ensures the exact estimation of the time-varying depth and mitigates singularity in the estimated rotation matrix. The proposed method eliminates the need to calibrate camera extrinsic parameters and measure the depths of reference feature points. Experiments on a 7-DoF robot manipulator have verified the effectiveness of the proposed CL-HBVS method.

## I. INTRODUCTION

Visual servoing, typically including position-based visual servoing (PBVS) and image-based visual servoing (IBVS), utilizes visual feedback to control robot motion [1]. PBVS relies on known scene geometry to reconstruct the six-degree-of-freedom (6-DoF) pose (i.e., position and orientation) from two-dimensional (2D) images, which is then used for robot control within the feedback loop. IBVS directly employs 2D images in the feedback loop with no prior knowledge about scenes but has several drawbacks, such as local minima and image Jacobian singularity. Homography-based visual servoing (HBVS) combines PBVS and IBVS to reconstruct a 6-DoF pose using only a pair of images, which greatly alleviates the drawbacks of both PBVS and IBVS [2].

Parametric uncertainties are common in vision-based robotic systems, where the time-varying depth of feature points and uncalibrated camera parameters are two typical types of parametric uncertainties [3]. Different adaptive control strategies have been proposed to handle these parameter uncertainties [4]–[9]. IBVS methods with passivity-based adaptive laws were developed for robot manipulators in [4]–[6], but they solely achieve pixel error convergence without considering 6-DoF pose control. Note that in the above methods, the time-varying

depth of image features can be expressed as a function of unknown parameters such that it can be accurately estimated if parameter estimates converge to their true values. Despite considering the time-varying depth, depth error convergence can not be achieved in these adaptive IBVS methods owing to the requirement of a stringent condition known as persistent excitation (PE) for parameter convergence. Adaptive HBVS methods were developed for robot manipulators in [7]–[9], but they require the reference extrinsic parameters of the camera to be precisely known a priori. Besides, they often assume a constant depth and rarely consider a time-varying depth.

This paper considers robotic manipulators under eye-to-hand (ETH) monocular cameras with unknown extrinsic parameters and proposes a passivity-based kinematic control solution named composite learning HBVS (CL-HBVS) to achieve 6-DoF pose tracking. First, the rotation matrix and depth ratio are extracted from homography decomposition; second, a linearly parameterized camera model is built by extracting extrinsic parameters and representing the time-varying depth via these parameters; third, an HBVS control law is proposed for 6-DoF pose tracking; lastly, a composite learning law is developed to estimate extrinsic parameters exactly, implying the exact estimation of the time-varying depth, under a condition of interval excitation (IE) that weakens PE [10]. Compared to existing adaptive HBVS methods, the distinctive feature of the proposed method is that it does not need to calibrate camera extrinsic parameters and measure the depths of reference feature points. Compared with our existing results on composite learning visual servoing [11]–[14], the current study has two distinctive features: 1) It is developed for visual tracking rather than visual regulation; 2) it is a kinematics-based design that facilitates industrial applications but takes the inner tracking error into consideration for performance enhancement.

Throughout this article,  $\mathbb{R}$ ,  $\mathbb{R}^+$ ,  $\mathbb{R}^n$ , and  $\mathbb{R}^{m \times n}$  are the spaces of real numbers, positive real numbers, real  $n$ -vectors, and real  $m \times n$ -matrices, respectively,  $\mathbb{N}$  is the set of natural numbers,  $\max\{\cdot\}$  is the maximum operator,  $\arccos(x)$  is the arc-cosine function,  $\det(A)$  is the determinant of  $A$ ,  $\|\mathbf{x}\|$  is the Euclidean norm,  $\ln(x)$  is the natural logarithm,  $\text{diag}(x_1, x_2, \dots, x_n)$  is a diagonal matrix with diagonal elements  $x_1$  to  $x_n$ , and  $x_i$  is the  $i$ th element of  $\mathbf{x}$  with  $i = 1$  to  $n$ , where  $x \in \mathbb{R}$ ,  $\mathbf{x} \in \mathbb{R}^n$ ,  $A \in \mathbb{R}^{n \times n}$ , and  $n, m, i \in \mathbb{N}$ .

## II. HOMOGRAPHY-BASED CAMERA MODEL

### A. Camera Projection and Euclidean Reconstruction

A robot manipulator for visual servoing tasks under an ETH monocular camera is depicted in Fig. 1. The setup comprises three Cartesian spatial frames: A world frame  $\{B\}$ , a camera

\*This work was supported in part by the Fundamental Research Funds for the Central Universities, Sun Yat-sen University, China, under Grant 23lgzy004 (Corresponding author: Yongping Pan).

<sup>1</sup>Yongping Pan is with the School of Advanced Manufacturing, Sun Yat-sen University, Shenzhen 518100, China panyongp@mail.sysu.edu.cn

<sup>2</sup>Beixian Lai and Zhiwen Li are with the School of Computer Science and Engineering, Sun Yat-sen University, Guangzhou 510006, China {laibx, lizhw63}@mail2.sysu.edu.cn

<sup>3</sup>Changyun Wen is with the School of Electrical and Electronic Engineering, Nanyang Technological University, Singapore 639798 ecywen@ntu.edu.sg

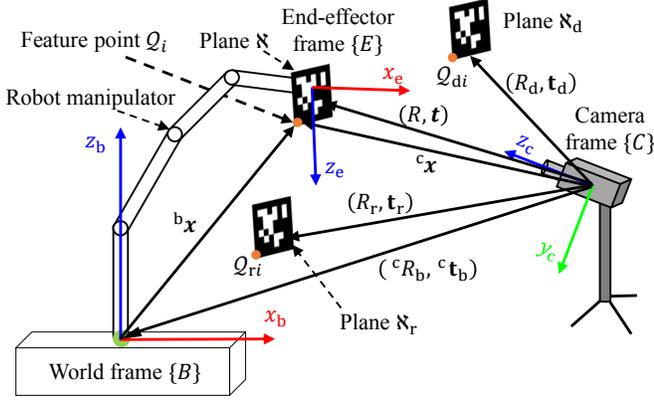


Fig. 1. A robot manipulator with an ETH monocular camera for HBVS.

frame  $\{C\}$ , and an end-effector frame  $\{E\}$ . It is assumed that: 1) The camera always captures a plane denoted as  $\aleph$ , which remains fixed in  $\{E\}$ ; 2) there exist  $N \geq 4$  ( $\in \mathbb{N}$ ) non-collinear feature points  $Q_i$  ( $i = 1$  to  $N$ ) lying on  $\aleph$ ; 3) the origin of  $\{E\}$  corresponds to the location of  $Q_1$ ; 4) a fixed and constant reference plane, denoted as  $\aleph_r$ , contains  $N$  reference feature points  $Q_{ri}$ ; 5) a time-varying desired plane named  $\aleph_d$  consists of  $N$  desired feature points  $Q_{di}$ .

Let  $R(t)$ ,  $R_r$ ,  $R_d(t) \in \mathbb{R}^{3 \times 3}$  denote the current, reference and desired orientations corresponding to  $\aleph$ ,  $\aleph_r$  and  $\aleph_d$  in  $\{C\}$ , respectively, and  $\mathbf{t}(t)$ ,  $\mathbf{t}_r$ ,  $\mathbf{t}_d(t) \in \mathbb{R}^3$  denote their positions in  $\{C\}$ , respectively. Then, let the Euclidean coordinates of  $Q_i$ ,  $Q_{ri}$  and  $Q_{di}$  in  $\{C\}$  be  ${}^c\mathbf{x}_i(t) = [x_i(t), y_i(t), z_i(t)]^T \in \mathbb{R}^3$ ,  ${}^c\mathbf{x}_{ri} = [x_{ri}, y_{ri}, z_{ri}]^T \in \mathbb{R}^3$  and  ${}^c\mathbf{x}_{di}(t) = [x_{di}(t), y_{di}(t), z_{di}(t)]^T \in \mathbb{R}^3$ , and their normalized coordinates be

$$\mathbf{x}_i(t) := {}^c\mathbf{x}_i/z_i, \quad \mathbf{x}_{ri} := {}^c\mathbf{x}_{ri}/z_{ri}, \quad \mathbf{x}_{di}(t) := {}^c\mathbf{x}_{di}/z_{di}$$

respectively. From the Euclidean geometry in Fig. 1, one gets the following relationships among  $\mathbf{x}_i$ ,  $\mathbf{x}_{di}$  and  $\mathbf{x}_{ri}$  [2]:

$$\mathbf{x}_i = \underbrace{\alpha_i}_{z_{ri}/z_i} \underbrace{H}_{\tilde{R} + (\tilde{\mathbf{t}}/d_r)\mathbf{n}_r^T} \mathbf{x}_{ri}, \quad (1)$$

$$\mathbf{x}_{di} = \underbrace{\alpha_{di}}_{z_{ri}/z_{di}} \underbrace{H_d}_{\tilde{R}_d + (\tilde{\mathbf{t}}_d/d_r)\mathbf{n}_r^T} \mathbf{x}_{ri} \quad (2)$$

in which  $\alpha_i(t) \in \mathbb{R}^+$  and  $\alpha_{di}(t) \in \mathbb{R}^+$  are the current and desired depth ratios, respectively,  $H(t) \in \mathbb{R}^{3 \times 3}$  and  $H_d(t) \in \mathbb{R}^{3 \times 3}$  are the current and desired Euclidean homography matrices, respectively,  $\tilde{R}(t) := R(t)R_r^T \in \mathbb{R}^{3 \times 3}$  and  $\tilde{\mathbf{t}}(t) := \mathbf{t}(t) - \tilde{R}(t)\mathbf{t}_r \in \mathbb{R}^3$  are the current mismatch rotation matrix and translation vector between  $\aleph$  and  $\aleph_r$ , respectively,  $\tilde{R}_d(t) := R_d(t)R_r^T \in \mathbb{R}^{3 \times 3}$  and  $\tilde{\mathbf{t}}_d(t) := \mathbf{t}_d(t) - \tilde{R}_d(t)\mathbf{t}_r \in \mathbb{R}^3$  are the desired mismatch rotation matrix and translation vector between  $\aleph_d$  and  $\aleph_r$ , respectively,  $\mathbf{n}_r \in \mathbb{R}^3$  is a constant unit normal of  $\aleph_r$  in  $\{C\}$ , and  $d_r \in \mathbb{R}^+$  is a constant distance between  $\{C\}$  and  $\aleph_r$  along  $\mathbf{n}_r$ . Let the homogeneous pixel coordinates of  $Q_i$ ,  $Q_{ri}$  and  $Q_{di}$  in the image plane be  $\mathbf{p}_i(t) = [u_i(t), v_i(t), 1]^T \in \mathbb{R}^3$ ,  $\mathbf{p}_{ri} = [u_{ri}, v_{ri}, 1]^T \in \mathbb{R}^3$  and  $\mathbf{p}_{di}(t) = [u_{di}(t), v_{di}(t), 1]^T \in \mathbb{R}^3$ , respectively. Using the perspective projection model of pinhole cameras, one obtains [15]

$$\mathbf{p}_i = K\mathbf{x}_i, \quad \mathbf{p}_{ri} = K\mathbf{x}_{ri}, \quad \mathbf{p}_{di} = K\mathbf{x}_{di} \quad (3)$$

where  $K \in \mathbb{R}^{3 \times 3}$  is an intrinsic parameter matrix. Combining (3) with (1) and (2), one obtains

$$\mathbf{p}_i = \alpha_i \underbrace{KHK^{-1}}_G \mathbf{p}_{ri}, \quad (4)$$

$$\mathbf{p}_{di} = \alpha_{di} \underbrace{KH_dK^{-1}}_{G_d} \mathbf{p}_{ri} \quad (5)$$

where  $G(t)$ ,  $G_d(t) \in \mathbb{R}^{3 \times 3}$  denote the current and desired projective homography matrices, respectively. The current mismatch rotation matrix  $\tilde{R}$  and depth ratio  $\alpha_i$  can be obtained from  $G$  by homography decomposition [14]. In addition, the desired ones  $\tilde{R}_d$  and  $\alpha_{di}$  can be obtained from  $G_d$ .

### B. Linear Parameterization of Camera Models

Let  ${}^cR_b \in \mathbb{R}^{3 \times 3}$  and  ${}^c\mathbf{t}_b \in \mathbb{R}^3$  denote a constant rotation matrix and a translation vector, respectively, which describe the transformation between the frames  $\{C\}$  and  $\{B\}$ . Let  ${}^b\mathbf{x}_i(\mathbf{q}) \in \mathbb{R}^4$  be a homogeneous Cartesian coordinate of the feature point  $Q_i$  in  $\{B\}$ , depending on the joint position  $\mathbf{q}(t) \in \mathbb{R}^n$  based on the robot forward kinematics [16], where  $n \in \mathbb{N}$  is the number of DoFs. Then, the perspective projection relationship in (3) can be rewritten as follows [14], [17]:

$$\mathbf{p}_i = KD^b\mathbf{x}_i(\mathbf{q})/z_i(\mathbf{q}) \quad (6)$$

in which  $D := [{}^cR_b, {}^c\mathbf{t}_b] \in \mathbb{R}^{3 \times 4}$  is an extrinsic parameter matrix, and  $z_i(\mathbf{q}) \in \mathbb{R}^+$  is the depth of  $Q_i$  in  $\{C\}$  given by

$$z_i(\mathbf{q}) = \mathbf{d}_3^T \mathbf{x}_i(\mathbf{q}) \quad (7)$$

with  $\mathbf{d}_3^T$  being the 3rd row of  $D$  [6]. Multiplying each side of (6) by  $z_i(\mathbf{q})$  and substituting (7) into the result yield

$$\mathbf{p}_i \mathbf{d}_3^T \mathbf{x}_i(\mathbf{q}) - KD^b\mathbf{x}_i(\mathbf{q}) = \mathbf{0}. \quad (8)$$

The following key property based on (8) from [18] is useful for the adaptive control design.

*Property 1:* Let  $\boldsymbol{\theta} := [d_{11}, d_{12}, \dots, d_{32}, d_{33}]^T \in \mathbb{R}^{11}$  be a camera's extrinsic parameter vector that is unknown but constant, where  $d_{ij} \in \mathbb{R}$  is the  $ij$ th element of  $D$ . If  $d_{34} \in \mathbb{R}^+$  is known, for any given  $\boldsymbol{\eta} \in \mathbb{R}^3$  and  $\gamma \in \mathbb{R}^4$ ,  $\boldsymbol{\eta} \mathbf{d}_3^T \gamma - KD\gamma$  can be linearly parameterized by

$$\boldsymbol{\eta} \mathbf{d}_3^T \gamma - KD\gamma = \Phi^T(\boldsymbol{\eta}, \gamma)\boldsymbol{\theta} - \mathbf{y}(\boldsymbol{\eta}, \gamma) \quad (9)$$

with  $\Phi(\boldsymbol{\eta}, \gamma) := B(\gamma)(\mathbf{k}\boldsymbol{\eta}^T - K^T) \in \mathbb{R}^{11 \times 3}$ ,  $\mathbf{y}(\boldsymbol{\eta}, \gamma) := d_{34}\gamma_4(K\mathbf{k} - \boldsymbol{\eta}) \in \mathbb{R}^3$ ,  $\mathbf{k} = [0, 0, 1]^T \in \mathbb{R}^3$ , and

$$B(\gamma) := \begin{bmatrix} \gamma & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \gamma & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \gamma_{[1:3]} \end{bmatrix} \in \mathbb{R}^{11 \times 3} \quad (10)$$

where  $\gamma_{[1:3]} \in \mathbb{R}^3$  contains the first 3 elements of  $\gamma$ . Then, substituting  $\boldsymbol{\eta} = \mathbf{p}_i$  and  $\gamma = {}^b\mathbf{x}_i$  into (9), one rewrites (8) into a linearly parameterized camera model

$$\mathbf{y}(\mathbf{p}_i, {}^b\mathbf{x}_i) = \Phi^T(\mathbf{p}_i, {}^b\mathbf{x}_i)\boldsymbol{\theta}, \quad i = 1 \text{ to } N. \quad (11)$$

### III. COMPOSITE LEARNING VISUAL SERVOING

It is assumed that the point  $\mathcal{Q}_1$ , denoted by  $\mathcal{Q}_e$ , is the origin of the end-effector frame  $\{E\}$ . Therefore,  $\mathcal{Q}_1$ ,  $\mathcal{Q}_{r1}$  and  $\mathcal{Q}_{d1}$  are replaced by  $\mathcal{Q}_e$ ,  $\mathcal{Q}_{re}$  and  $\mathcal{Q}_{de}$ , respectively. The control objective can be described by

$$\tilde{R} \rightarrow \tilde{R}_d, \alpha \rightarrow \alpha_d, \mathbf{p} \rightarrow \mathbf{p}_d \text{ as } t \rightarrow \infty. \quad (12)$$

With the consideration that the plane  $\aleph$  remains fixed in  $\{E\}$ , the orientation of the robot end-effector is represented by  $R$ . Consequently, the triplet  $(R, z, \mathbf{p})$  uniquely determines the pose of the end-effector in the camera frame  $\{C\}$ . When  $\tilde{R} = \tilde{R}_d$  and  $\alpha = \alpha_d$ , one obtains  $R = R_d$  and  $z = z_d$ . Utilizing  $z\mathbf{p} = K^c \mathbf{x}$ , if  $\mathbf{p} = \mathbf{p}_d$ , one obtains  ${}^c \mathbf{x} = {}^c \mathbf{x}_d$ , indicating that the end-effector reaches its desired pose in  $\{C\}$ .

#### A. Pose Error Definition

For control synthesis, definite a rotation error [2]

$$\mathbf{e}_\omega(t) := \boldsymbol{\vartheta}(t) - \boldsymbol{\vartheta}_d(t) \quad (13)$$

where  $\boldsymbol{\vartheta}(t), \boldsymbol{\vartheta}_d(t) \in \mathbb{R}^3$  are the angle-axis representations of  $\tilde{R}$  and  $\tilde{R}_d$ , respectively, given by

$$\boldsymbol{\vartheta}(t) := \boldsymbol{\mu}(t)\phi(t), \boldsymbol{\vartheta}_d(t) := \boldsymbol{\mu}_d(t)\phi_d(t) \quad (14)$$

in which  $\phi(t), \phi_d(t) \in \mathbb{R}$  are the rotation angles around unit axes  $\boldsymbol{\mu}(t), \boldsymbol{\mu}_d(t) \in \mathbb{R}^3$ , respectively, confined to

$$-\pi < \phi(t) < \pi, -\pi < \phi_d(t) < \pi.$$

The solutions for  $\phi$  and  $\boldsymbol{\mu}$  can be determined by

$$\begin{cases} \phi(t) = \arccos(\frac{1}{2}(\text{tr}(\tilde{R}) - 1)), \\ [\boldsymbol{\mu}(t)]_\times = \frac{\tilde{R} - \tilde{R}^T}{2 \sin(\phi)} \end{cases} \quad (15)$$

where  $[\cdot]_\times$  is a skew symmetry operator [14]. Likewise, with  $\tilde{R}_d$  being obtained from homography decomposition, the solutions for  $\phi_d$  and  $\boldsymbol{\mu}_d$  can also be determined. Differentiating  $\mathbf{e}_\omega$  with respect to time  $t$  yields

$$\dot{\mathbf{e}}_\omega(t) = L_\omega(\boldsymbol{\vartheta})^c R_b {}^b \boldsymbol{\omega} - \dot{\boldsymbol{\vartheta}}_d \quad (16)$$

in which  ${}^b \boldsymbol{\omega}(t) \in \mathbb{R}^3$  is an angular velocity of  $\mathcal{Q}_e$  in  $\{B\}$ , and  $L_\omega(\boldsymbol{\vartheta}) \in \mathbb{R}^{3 \times 3}$  is a Jacobian-like matrix given by

$$L_\omega(\boldsymbol{\vartheta}) := I - \frac{\phi}{2} [\boldsymbol{\mu}]_\times + \left(1 - \frac{\text{sinc}(\phi)}{\text{sinc}^2(\phi/2)}\right) [\boldsymbol{\mu}]_\times^2$$

with  $\text{sinc}(\phi) := \sin(\phi)/\phi$  and  $\text{sinc}(0) = 1$ .

To control the translation of the feature point  $\mathcal{Q}_e$ , define current and desired extended translation vectors

$$\begin{cases} \mathbf{p}_e(t) := [u(t), v(t), -\ln(\alpha(t))]^T, \\ \mathbf{p}_{ed}(t) := [u_d(t), v_d(t), -\ln(\alpha_d(t))]^T \end{cases} \quad (17)$$

respectively. Then, define an extended translation error

$$\mathbf{e}_v(t) := \mathbf{p}_e(t) - \mathbf{p}_{ed}(t). \quad (18)$$

Differentiating  $z(\mathbf{q})$  in (7) with respect to time  $t$  yields

$$\dot{z}(\mathbf{q}) = \mathbf{d}_3^T {}^b \mathbf{v}, \quad (19)$$

where  ${}^b \mathbf{v} := {}^b \dot{\mathbf{x}}(\mathbf{q}) \in \mathbb{R}^4$  is a linear velocity of  $\mathcal{Q}_e$  in  $\{B\}$ . From (6), (7) and (19), the time derivative of  $\mathbf{e}_v$  becomes

$$\dot{\mathbf{e}}_v = A_e(\mathbf{p}) {}^b \mathbf{v} / z(\mathbf{q}) - \dot{\mathbf{p}}_{ed} \quad (20)$$

in which  $A_e(\mathbf{p}) := KD - \mathbf{p}_0 \mathbf{d}_3^T \in \mathbb{R}^{3 \times 4}$  is an extended interaction matrix with  $\mathbf{p}_0(t) := [u(t), v(t), 0]^T \in \mathbb{R}^3$ .

Based on the robot forward kinematics [16], the velocities  ${}^b \mathbf{v}$  and  ${}^b \boldsymbol{\omega}$  are calculated by

$${}^b \mathbf{v} = J_v(\mathbf{q}) \dot{\mathbf{q}}, {}^b \boldsymbol{\omega} = J_\omega(\mathbf{q}) \dot{\mathbf{q}} \quad (21)$$

where  $J_v(\mathbf{q}) \in \mathbb{R}^{4 \times n}$  and  $J_\omega(\mathbf{q}) \in \mathbb{R}^{3 \times n}$  are the translational and rotational Jacobian matrices, respectively, which can be obtained from the robot's forward kinematics. Inspired by [5], multiplying each side of (20) by  $z(\mathbf{q})$  and substituting (19) and (21) into the obtained result and (20), one obtains an overall kinematic system for 6-DoF pose control as follows:

$$\begin{bmatrix} z(\mathbf{q}) \dot{\mathbf{e}}_v + \frac{1}{2} \dot{z}(\mathbf{q}) \mathbf{e}_v \\ \dot{\mathbf{e}}_\omega \end{bmatrix} = J_p(\mathbf{p}, \mathbf{q}, \mathbf{e}_v, \boldsymbol{\vartheta}) \dot{\mathbf{q}} - \begin{bmatrix} z(\mathbf{q}) \dot{\mathbf{p}}_{ed} \\ \dot{\boldsymbol{\vartheta}}_d \end{bmatrix} \quad (22)$$

with  $J_p(\mathbf{p}, \mathbf{q}, \mathbf{e}_v, \boldsymbol{\vartheta}) := [J_{pv}^T(\mathbf{p}, \mathbf{q}, \mathbf{e}_v), J_{p\omega}^T(\mathbf{q}, \boldsymbol{\vartheta})]^T \in \mathbb{R}^{6 \times n}$  being a Jacobian matrix that maps the joint velocity  $\dot{\mathbf{q}}$  to the velocity pairs  $(z(\mathbf{q}) \dot{\mathbf{p}}_e + \frac{1}{2} \dot{z}(\mathbf{q}) \mathbf{e}_v, \dot{\boldsymbol{\vartheta}})$ . Here,  $J_{pv}(\mathbf{p}, \mathbf{q}, \mathbf{e}_v) := Q(\mathbf{p}, \mathbf{e}_v) J_v(\mathbf{q}) \in \mathbb{R}^{3 \times n}$  with  $Q(\mathbf{p}, \mathbf{e}_v) := A_e(\mathbf{p}) + \frac{1}{2} \mathbf{e}_v \mathbf{d}_3^T \in \mathbb{R}^{3 \times 4}$  is the translational mapping component, and  $J_{p\omega}(\mathbf{q}, \boldsymbol{\vartheta}) := L_\omega(\boldsymbol{\vartheta})^c R_b J_\omega(\mathbf{q}) \in \mathbb{R}^{3 \times n}$  is the rotational mapping component. To simplify notation, define two auxiliary variables  $\mathbf{u}_v := J_{pv}(\mathbf{p}, \mathbf{q}, \mathbf{e}_v) \dot{\mathbf{q}} - z(\mathbf{q}) \dot{\mathbf{p}}_{ed} \in \mathbb{R}^3$  and  $\mathbf{u}_\omega := J_{p\omega}(\mathbf{q}, \boldsymbol{\vartheta}) \dot{\mathbf{q}} - \dot{\boldsymbol{\vartheta}}_d \in \mathbb{R}^3$ . Then, it follows from [19] that (22) is passive concerning input-output pairs  $(\mathbf{u}_v, \mathbf{e}_v)$  and  $(\mathbf{u}_\omega, \mathbf{e}_\omega)$  for the translational and rotational components, respectively, and the corresponding storage functions are  $V_v = z(\mathbf{q}) \mathbf{e}_v^T \mathbf{e}_v / 2$  and  $V_\omega = \mathbf{e}_\omega^T \mathbf{e}_\omega / 2$ , respectively. With the translational and rotational errors  $\mathbf{e}_\omega$  and  $\mathbf{e}_v$ , the control objective (12) becomes

$$\mathbf{e}_\omega \rightarrow \mathbf{0}, \mathbf{e}_v \rightarrow \mathbf{0} \text{ as } t \rightarrow \infty. \quad (23)$$

#### B. Homography-Base Visual Servoing

Let  $(\mathbf{p}_{ed}, \dot{\mathbf{p}}_{ed}, \ddot{\mathbf{p}}_{ed})$  and  $(\boldsymbol{\vartheta}_d, \dot{\boldsymbol{\vartheta}}_d, \ddot{\boldsymbol{\vartheta}}_d)$  be desired extended translational and rotational trajectories of the feature point  $\mathcal{Q}_{ed}$ , respectively. Then, introduce nominal reference translational and rotational trajectories as follows:

$$\mathbf{p}_{er} := \dot{\mathbf{p}}_{ed} - \lambda_v \mathbf{e}_v, \boldsymbol{\vartheta}_r := \dot{\boldsymbol{\vartheta}}_d - \lambda_\omega \mathbf{e}_\omega \quad (24)$$

respectively, where  $\dot{\mathbf{p}}_{er}, \dot{\boldsymbol{\vartheta}}_r \in \mathbb{R}^3$  denote the reference translation and rotation velocities, respectively, and  $\lambda_v, \lambda_\omega \in \mathbb{R}^+$  are certain constants. Let  $\hat{\boldsymbol{\theta}} \in \mathbb{R}^{11}$  be an estimate of  $\boldsymbol{\theta}$ . Substituting  $\hat{\boldsymbol{\theta}}$  into  $J_p(\mathbf{p}, \mathbf{q}, \boldsymbol{\vartheta}, \mathbf{e}_v)$ ,  $Q(\mathbf{p}, \mathbf{e}_v)$ ,  ${}^c R_b$ , and  $\mathbf{d}_3$ , one gets

$$\begin{aligned} \hat{J}_p &:= J_p|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}}, \hat{Q} := Q|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}}, {}^c \hat{R}_b := {}^c R_b|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}}, \hat{z} := z|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}} \\ \text{with } \hat{J}_p(\mathbf{p}, \mathbf{q}, \boldsymbol{\vartheta}, \mathbf{e}_v) &\in \mathbb{R}^{6 \times n}, \hat{Q}(\mathbf{p}, \mathbf{e}_v) \in \mathbb{R}^{3 \times 4}, {}^c \hat{R}_b \in \mathbb{R}^{3 \times 3} \\ \text{and } \hat{z} &\in \mathbb{R}. \text{ Then, define estimation errors} \\ \tilde{J}_p &:= \hat{J}_p - J_p, \tilde{Q} := \hat{Q} - Q, {}^c \tilde{R}_b := {}^c \hat{R}_b - {}^c R_b, \tilde{z} := \hat{z} - z. \end{aligned}$$

Design an HBVS-based joint velocity control law

$$\begin{aligned} \dot{\mathbf{q}}_c &= \hat{J}_p^+(\mathbf{p}, \mathbf{q}, \boldsymbol{\vartheta}, \mathbf{e}_v) \begin{bmatrix} \hat{z} \dot{\mathbf{p}}_{er} \\ \dot{\boldsymbol{\vartheta}}_r \end{bmatrix} - J_v^T(\mathbf{q}) \hat{Q}^T(\mathbf{p}, \mathbf{e}_v) K_p \mathbf{e}_v \\ &\quad - J_\omega^T(\mathbf{q}) {}^c \hat{R}_b^T L_\omega^T(\boldsymbol{\vartheta}) K_\omega \mathbf{e}_\omega \end{aligned} \quad (25)$$

where  $K_p, K_\omega \in \mathbb{R}^{3 \times 3}$  are positive-definite diagonal matrices of control gains, and  $\hat{J}_p^+(\mathbf{p}, \mathbf{q}, \boldsymbol{\vartheta}, \mathbf{e}_v) \in \mathbb{R}^{n \times 6}$  is the pseudo inverse of  $\hat{J}_p(\mathbf{p}, \mathbf{q}, \boldsymbol{\vartheta}, \mathbf{e}_v)$ . It is observed from (25) that the estimated matrices  $\hat{J}_p^+, \hat{Q}$  and  ${}^c\hat{R}_b$  determine how the errors  $\mathbf{e}_v, \mathbf{e}_\omega$ , and velocities  $\dot{\mathbf{p}}_{er}$  and  $\dot{\boldsymbol{\vartheta}}_r$  are projected into the joint space, which affects the performance of visual servoing. Then, define a velocity tracking error in the joint space

$$\dot{\mathbf{e}}(t) := \dot{\mathbf{q}}(t) - \dot{\mathbf{q}}_c(t) \quad (26)$$

and two auxiliary matrices  $\Phi_{y_i} := \Phi(\mathbf{p}_0 - \mathbf{e}_v/2, \mathbf{J}_{v_i}) \in \mathbb{R}^{11 \times 3}$  and  $\Phi_{r_i} := B(\gamma)|_{\gamma=[J_{\omega_i}^T, 0]^T} \in \mathbb{R}^{11 \times 3}$ , where  $\mathbf{J}_{v_i} \in \mathbb{R}^4$  and  $\mathbf{J}_{\omega_i} \in \mathbb{R}^3$  are the  $i$ th columns of  $\mathbf{J}_v(\mathbf{q})$  and  $\mathbf{J}_\omega(\mathbf{q})$ , respectively. Then, one has two linearly parameterized models.

*Property 2:* If  $d_{34}$  is known, then for any given  $\boldsymbol{\eta}_1$ , one gets a linearly parameterized model

$$\mathbf{J}_v^T(\mathbf{q})\hat{Q}^T(\mathbf{p}, \mathbf{e}_v)\boldsymbol{\eta}_1 + \mathbf{J}_\omega^T(\mathbf{q}){}^c\hat{R}_b^T\boldsymbol{\eta}_2 = \mathbf{Y}_1^T(\mathbf{p}, \mathbf{q}, \boldsymbol{\eta}_1, \boldsymbol{\eta}_2)\hat{\boldsymbol{\theta}} \quad (27)$$

with  $\mathbf{Y}_1(\mathbf{p}, \mathbf{q}, \boldsymbol{\eta}_1, \boldsymbol{\eta}_2) := [\Phi_{r1}\boldsymbol{\eta}_2 - \Phi_{y1}\boldsymbol{\eta}_1, \Phi_{r2}\boldsymbol{\eta}_2 - \Phi_{y2}\boldsymbol{\eta}_1, \dots, \Phi_{rn}\boldsymbol{\eta}_2 - \Phi_{yn}\boldsymbol{\eta}_1] \in \mathbb{R}^{11 \times n}$ .

*Property 3:* If  $d_{34}$  is known, then for any given  $\boldsymbol{\eta} \in \mathbb{R}^3$ , one gets a linearly parameterized model

$$\hat{\mathbf{z}}\boldsymbol{\eta} = \mathbf{Y}_2^T(\mathbf{q}, \boldsymbol{\eta})\hat{\boldsymbol{\theta}} + \mathbf{y}_2(\boldsymbol{\eta}) \quad (28)$$

with  $\mathbf{y}_2(\boldsymbol{\eta}) := d_{34}\boldsymbol{\eta} \in \mathbb{R}^3$ ,  $\mathbf{Y}_2(\mathbf{q}, \boldsymbol{\eta}) := B({}^b\mathbf{x}(\mathbf{q}))\mathbf{k}\boldsymbol{\eta}^T \in \mathbb{R}^{11 \times 3}$ , and  $B(\cdot)$  defined in (10).

Replace the tuple  $(\hat{\mathbf{z}}, \hat{\boldsymbol{\theta}}, \boldsymbol{\eta})$  in (28) by  $(z(\mathbf{q}), \boldsymbol{\theta}, \dot{\mathbf{p}}_{er})$  yields

$$(\hat{\mathbf{z}} - z(\mathbf{q}))\dot{\mathbf{p}}_{er} = \mathbf{Y}_2^T(\mathbf{q}, \dot{\mathbf{p}}_{er})\tilde{\boldsymbol{\theta}} \quad (29)$$

where  $\tilde{\boldsymbol{\theta}}(t) := \hat{\boldsymbol{\theta}}(t) - \boldsymbol{\theta} \in \mathbb{R}^{11}$  is a parameter estimation error. By using (27), (25) can be rewritten into

$$\dot{\mathbf{q}}_c = \hat{J}_p^+(\mathbf{p}, \mathbf{q}, \boldsymbol{\vartheta}, \mathbf{e}_v) \begin{bmatrix} \hat{\mathbf{z}}\dot{\mathbf{p}}_{er} \\ \dot{\boldsymbol{\vartheta}}_r \end{bmatrix} - \mathbf{Y}_1^T(\mathbf{p}, \mathbf{q}, \boldsymbol{\eta}_1, \boldsymbol{\eta}_2)\hat{\boldsymbol{\theta}} \quad (30)$$

with  $\boldsymbol{\eta}_1 = K_p\mathbf{e}_v$  and  $\boldsymbol{\eta}_2 = L_\omega^T K_\omega \mathbf{e}_\omega$ . Applying (24), (26) and (30) to (22), yields the closed-loop overall kinematics

$$\begin{bmatrix} z\dot{\mathbf{e}}_v + \frac{1}{2}\dot{\mathbf{z}}\mathbf{e}_v \\ \dot{\mathbf{e}}_\omega \end{bmatrix} = \hat{J}_p\dot{\mathbf{e}} - \hat{J}_p\mathbf{Y}_1^T\hat{\boldsymbol{\theta}} - \tilde{J}_p\dot{\mathbf{e}} + \tilde{J}_p\mathbf{Y}_1^T\hat{\boldsymbol{\theta}} - \tilde{J}_p\hat{J}_p^+ \begin{bmatrix} \hat{\mathbf{z}}\dot{\mathbf{p}}_{er} \\ \dot{\boldsymbol{\vartheta}}_r \end{bmatrix} + \begin{bmatrix} \tilde{\mathbf{z}}\dot{\mathbf{p}}_{er} \\ \mathbf{0} \end{bmatrix} - \begin{bmatrix} \lambda_v z \mathbf{e}_v \\ \lambda_\omega \mathbf{e}_\omega \end{bmatrix}. \quad (31)$$

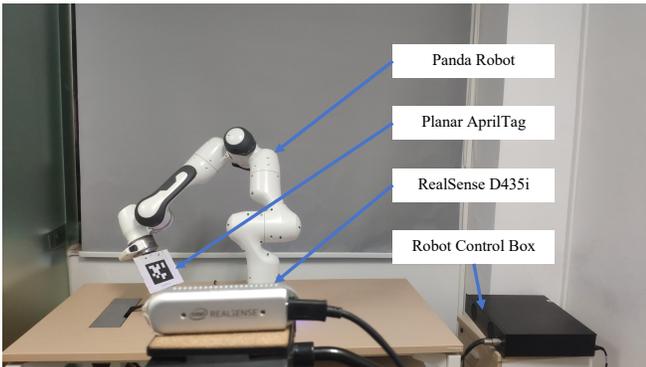


Fig. 2. An experimental environment for monocular ETH robots, where the end-effector is identified by an attached planar AprilTag.

### C. Composite Learning Adaptation

We define IE and PE to facilitate analysis as follows [10].

*Definition 1:* A bounded signal  $\Phi(t) \in \mathbb{R}^{11 \times 3}$  is of IE if  $\exists T_e, \zeta_d, \sigma \in \mathbb{R}^+$  such that  $\int_{T_e - \zeta_d}^{T_e} \Phi(\zeta)\Phi^T(\zeta)d\zeta \geq \sigma I$ .

*Definition 2:* A bounded signal  $\Phi(t) \in \mathbb{R}^{11 \times 3}$  is of PE if  $\exists \zeta_d, \sigma \in \mathbb{R}^+$  such that  $\int_{t - \zeta_d}^t \Phi(\zeta)\Phi^T(\zeta)d\zeta \geq \sigma I, \forall t \geq 0$ .

Define a generalized prediction error

$$\boldsymbol{\xi}(t) := \Psi_e(t)\hat{\boldsymbol{\theta}} - \Psi_e(t)\boldsymbol{\theta} \in \mathbb{R}^{11} \quad (32)$$

in which  $\Psi_e \in \mathbb{R}^{11 \times 11}$  is given by

$$\Psi_e(t) := \begin{cases} \int_{t - \zeta_d}^t \Phi(\zeta)\Phi^T(\zeta)d\zeta, & t < T_e \\ \int_{T_e - \zeta_d}^{T_e} \Phi(\zeta)\Phi^T(\zeta)d\zeta, & t \geq T_e \end{cases}.$$

Note that one has  $\Psi_e\boldsymbol{\theta} = \int_{t - \zeta_d}^t \Phi(\zeta)\boldsymbol{\theta}d\zeta$ . Using eigendecomposition, one obtains  $\Psi_e = U\Sigma U^T$  with  $\Sigma := \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_{11})$ , where  $\lambda_i \in \mathbb{R}$  ( $i = 1$  to  $11$ ) is the  $i$ th eigenvalue of  $\Psi_e$ , and  $U \in \mathbb{R}^{11 \times 11}$  is orthogonal, i.e.,  $UU^T = I$ . Note that  $\Psi_e$  is symmetric and positive-semidefinite, and each eigenvalue  $\lambda_i$  is non-negative. Define a regularized inversion of  $\Psi_e$ :

$$\Psi_\rho^+ := U\Sigma_\rho U^T \in \mathbb{R}^{11 \times 11} \quad (33)$$

with  $\Sigma_\rho := \text{diag}(\lambda_{\rho 1}, \lambda_{\rho 2}, \dots, \lambda_{\rho 11}) \in \mathbb{R}^{11 \times 11}$ , where  $\lambda_{\rho i} := 1/\max\{\lambda_i, \rho\} \in \mathbb{R}^+$  ( $i = 1$  to  $11$ ) is the  $i$ th eigenvalue of  $\Psi_\rho^+$ , and  $\rho \in \mathbb{R}^+$  is a small threshold. Then, one concludes that: 1) If  $\Psi_e \in \mathbb{R}^{11 \times 11}$  is positive-semidefinite,  $\Psi_\rho^+$  in (33) is positive-definite; 2)  $\mathbf{0} \leq \Psi_\rho^+\Psi_e \leq I$ ,  $\Psi_\rho^+\Psi_e = I$  iff  $\Psi_e \geq \rho I$ , and  $\Psi_\rho^+\Psi_e = \mathbf{0}$  iff  $\Psi_e = \mathbf{0}$  [14]. Now, design a regularized composite learning law as follows:

$$\dot{\tilde{\boldsymbol{\theta}}} = \Gamma(\mathbf{Y}_1(\mathbf{p}, \mathbf{q}, \boldsymbol{\eta}_1, \boldsymbol{\eta}_2)\dot{\mathbf{q}} - \mathbf{Y}_2(\mathbf{q}, \boldsymbol{\eta})K_p\mathbf{e}_v - \kappa\Psi_\rho^+\boldsymbol{\xi}) \quad (34)$$

with  $\boldsymbol{\eta} = \dot{\mathbf{p}}_{er}$ , where  $\Gamma \in \mathbb{R}^{11 \times 11}$  is a positive-definite matrix of learning rates, and  $\kappa \in \mathbb{R}^+$  is a weighting factor. We give the following reasonable assumption from [5].

*Assumption 1:* The low-level joint controller of the robotic system guarantees  $\dot{\mathbf{e}} \in L_2 \cap L_\infty$  so that there exists a constant  $l_m \in \mathbb{R}^+$  to satisfy  $\int_0^t \dot{\mathbf{e}}^T(\tau)\dot{\mathbf{e}}(\tau)d\tau \leq l_m, \forall t \geq 0$ .

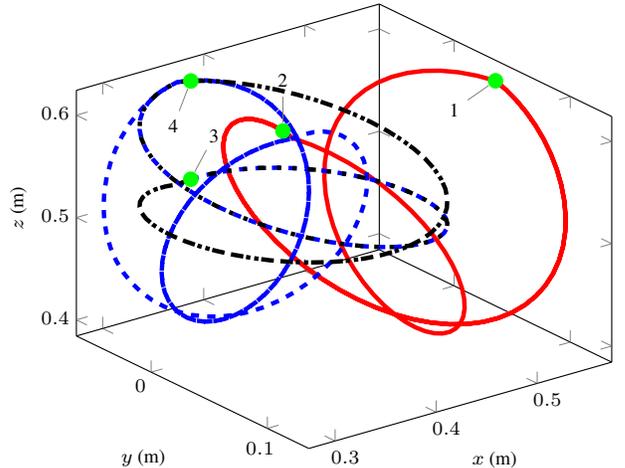


Fig. 3. The desired trajectory of the feature point  $Q_e$  in the Euclidean space, where the notation “1” to “4” (marked in green dot) denote the waypoints at instants  $t = 0, 30, 60$  and  $100$  s, respectively.

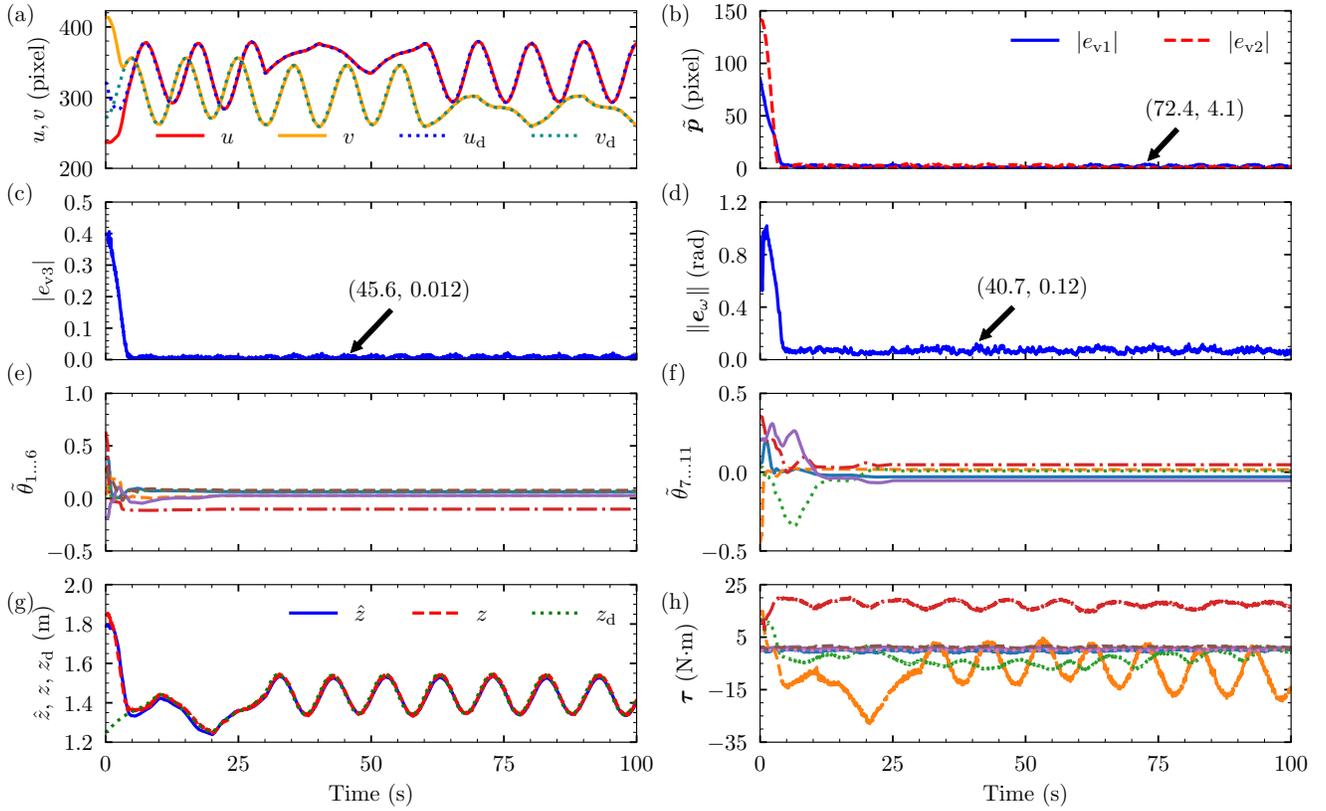


Fig. 4. Control results of the proposed CL-HBVS in experiments. (a) Pixel tracking trajectories in the  $u$ - and  $v$ -axes of the image plane. (b) Pixel errors  $e_{v1}$  and  $e_{v2}$  in the  $u$ - and  $v$ -axes of the image plane, respectively. (c) Absolute value of the depth ratio error  $e_{v3}$ . (d) Norm of the rotation error  $e_{\omega}$ . (e) Parameter errors  $\hat{\theta}_i$  ( $i = 1$  to  $6$ ). (f) Parameter errors  $\hat{\theta}_i$  ( $i = 7$  to  $11$ ). (g) Estimated, current, and desired depths  $\hat{z}$ ,  $z$ , and  $z_d$ . (h) Control torque  $\tau \in \mathbb{R}^7$ .

Consider an articulated robot system with an ETH monocular camera, as depicted in Fig. 1, where a plane  $\mathfrak{N}$  is attached to the end-effector frame  $\{E\}$ , containing at least 4 feature points. Only the feature point  $Q_e$  in the world frame  $\{B\}$  is known, which is applied to calculate its Cartesian coordinate  ${}^b\mathbf{x}(q)$ . If the kinematic system described by (22) is driven by the CL-HBVS control law (25) with (34), then the closed-loop system is stable in the following sense: 1) All signals involved are bounded and the tracking errors  $e_v(t)$  and  $e_{\omega}(t)$  asymptotically converge to  $\mathbf{0}$  on  $t \geq 0$ ; 2) if IE exists, both the tracking errors  $e_v(t)$ ,  $e_{\omega}(t)$  and the parameter estimation error  $\hat{\theta}(t)$  asymptotically converge to  $\mathbf{0}$  on  $t \geq T_e$ , implying exact depth estimation as  $\hat{z}(t) \rightarrow z(q(t))$ . The proof is similar to [5], so it is omitted here. Note that even if  $d_{34}$  is unknown a priori, (34) still works as discussed in [14].

#### IV. EXPERIMENTAL STUDIES

Experiments on the vision-based robotic system are carried out as illustrated in Fig. 2, where the environment includes a 7-DoF collaborative robot named Franka Emika Panda, and a fixed Intel RealSense camera D435i to capture images at a resolution of  $640 \times 480$  pixels, which provides visual signals at a rate of 30 frames per second. An AprilTag [20] with four visual markers positioned at its corners is affixed to the end-effector. During experiments, pixel positions  $p_{ri}$  ( $i = 1$  to  $4$ ) of each reference feature points  $Q_{ri}$  are (189, 241), (268, 241),

(266, 163) and (185, 163), and the initial robot joint position is  $q(0) = [0, -0.5236, -1.8326, -2.4435, -0.2618, 2.0071, -1.5708]^T$  rad. The camera's intrinsic parameter matrix is

$$K = \begin{bmatrix} 608 & 0 & 327 \\ 0 & 608 & 232 \\ 0 & 0 & 1 \end{bmatrix}.$$

Since the camera's extrinsic parameter  $D$  is unavailable in practice, the following calibrated one from ViSP [21]:

$$\hat{D}_r = \begin{bmatrix} -0.1184 & 0.9902 & -0.0735 & 0.1382 \\ 0.0427 & -0.0689 & -0.9967 & 0.6554 \\ -0.9920 & -0.1212 & -0.0341 & 1.8341 \end{bmatrix}$$

is chosen as a reference. Also, as  $d_{34}$  is unknown, a reference one  $\hat{d}_{34r} = 1.8341$  obtained from  $\hat{D}_r$  is used for calculating  $\mathbf{y}(p_i, {}^b\mathbf{x}_i)$  in (11) and  $\mathbf{y}_2(\eta)$  in (28).

An initial estimate of  $D$  is set as  $\hat{\theta}(0) = [0.2263, 0.9567, 0.1830, 1.0000, -0.1830, 0.2263, -0.9567, 0.1000, -0.9567, 0.1830, 0.2263]^T$ , which ensures a suitable initial discrepancy with  $D$ , better showing the evolution of the parameter estimate  $\hat{\theta}$ . We design a desired trajectory with 100 s shown in Fig. 3, which consists of three sub-tasks: 1) Tracking the red trajectory at  $t \in [0, 30)$  s; 2) tracking the blue trajectory at  $t \in [30, 60)$  s; 3) tracking the black trajectory at  $t \in [60, 100)$  s.

The proposed CL-HBVS (25) with (34) is implemented with  $K_p = \text{diag}(4 \times 10^{-5}, 4 \times 10^{-5}, 0.2)$ ,  $K_{\omega} = 0.2I$ ,  $\Gamma =$

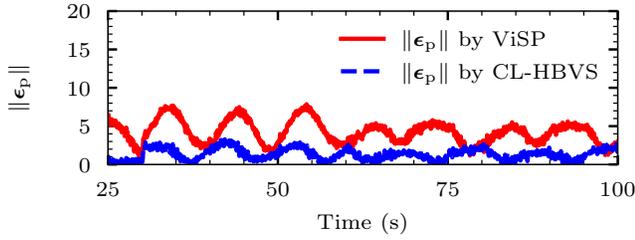


Fig. 5. Prediction errors  $\epsilon_p$  by ViSP and the proposed CL-HBVS.

TABLE I  
ESTIMATED, CURRENT AND DESIRED DEPTH VALUES  $\hat{z}$ ,  $z$ , AND  $z_d$   
DURING EXPERIMENTS

Time (s)	31	42	53	64	75	86	95	100
$\hat{z}$ (m)	1.48	1.51	1.53	1.51	1.47	1.40	1.46	1.40
$z$ (m)	1.49	1.52	1.54	1.52	1.48	1.41	1.47	1.42
$z_d$ (m)	1.50	1.53	1.54	1.51	1.46	1.40	1.45	1.43
$\frac{\ \hat{\theta}\ }{\ \theta\ }$ (%)	4.88	4.87	4.88	4.88	4.88	4.88	4.88	4.88

$0.4I$ ,  $\lambda_v = \lambda_\omega = 2$ ,  $\kappa = 15$ ,  $\varrho = 3$ , and  $T_e = \tau_d = 25$ . Experiments are run on a host PC under Linux OS Ubuntu 20.04 with ROS Noetic Ninjemys. The ViSP and the Franka Control Interface provided via libfranka are integrated into the ROS. To evaluate the effectiveness of depth tracking and estimation by the proposed method, we use the pose estimation algorithm embedded in ViSP to get exact information, where  $z$  is calculated by the known geometric information of the AprilTag, and  $\hat{z}$  is obtained by (7). Finally, define a model prediction error  $\epsilon_p(t) := \mathbf{y}(\mathbf{p}, {}^b\mathbf{x}) - \Phi^T(\mathbf{p}, {}^b\mathbf{x})\hat{\theta}$ , where  $\hat{\theta}$  is obtained by (34) ( $\kappa = 15$ ) and  $\hat{\theta} = \hat{\theta}_r$  for the proposed CL-HBVS and ViSP, respectively, and  $\hat{\theta}_r \in \mathbb{R}^{11}$  is the first 11 elements of the calibrated matrix  $\hat{D}_r$ .

Experiment results are demonstrated in Fig. 4. One observes that the current pixel  $\mathbf{p} = [u, v, 1]^T$  of the feature point  $Q_e$  accurately tracks its desired position  $\mathbf{p}_d = [u_d, v_d, 1]^T$  with a steady-state error of about 0.4 pixel [see Figs. 4(a) and (b)], the current depth  $z$  converges to the desired depth  $z_d$  with an error of about 0.02 m [see Fig. 4 (c)], and the norm of the rotation error  $\|e_\omega\|$  is about 0.1 rad on average [see Fig. 4(d)]. Besides, the convergence of  $\tilde{\theta}$  with an error of 0.1 is achieved at about 25 s [see Fig. 4 (e)-(f)], which indicates an accurate depth estimation [see Fig. 4 (g)]. Table I gives the depth values  $\hat{z}$ ,  $z$  and  $z_d$  at some instants, where both the depth estimation error  $\tilde{z}$  and the depth tracking error  $e_z := z - z_d \in \mathbb{R}$  are less than 2 cm. In addition, the control torque  $\tau \in \mathbb{R}^7$  [see Fig. 4 (h)] is shown to verify that the control input of the low-level joint controller is reasonable. What's more, the prediction error  $\epsilon_p$  obtained by the proposed CL-HBVS is smaller than that by ViSP [see Fig. 5], indicating that  $\hat{\theta}$  by the CL-HBVS is closer to the true value  $\theta$  compared to the calibrated one  $\hat{\theta}_r$ . To assess the performance of  $\tilde{\theta}$ , the percentage estimation error  $\|\tilde{\theta}\|/\|\theta_r\|$  (%) with less than 5 % is also shown in Table I.

## V. CONCLUSIONS

This paper has presented a passivity-based kinematic control method named CL-HBVS for robots with ETH monocular

cameras to achieve 6-DoF pose tracking in the 3D space under unknown camera extrinsic parameters. Parameter convergence is achieved under the weakened IE condition, which leads to the exact estimation of the time-varying depth and prevents singularity in the estimated rotation matrix. Experiments based on a 7-DoF robot manipulator have validated the effectiveness of the proposed method in estimating the extrinsic parameters and time-varying depth and controlling 6-DoF robot pose.

## REFERENCES

- [1] F. Chaumette and S. Hutchinson, "Visual servo control. I. basic approaches," *IEEE Robot. Autom. Mag.*, vol. 13, no. 4, pp. 82–90, Dec. 2006.
- [2] J. Chen, D. M. Dawson, W. E. Dixon, and A. Behal, "Adaptive homography-based visual servo tracking for a fixed camera configuration with a camera-in-hand extension," *IEEE Trans. Control Syst. Technol.*, vol. 13, no. 5, pp. 814–825, Sep. 2005.
- [3] C.-C. Cheah, C. Liu, and J.-J. E. Slotine, "Adaptive tracking control for robots with unknown kinematic and dynamic properties," *Int. J. Robot. Res.*, vol. 25, no. 3, pp. 283–296, Mar. 2006.
- [4] H. Wang, C. C. Cheah, W. Ren, and Y. Xie, "Passive separation approach to adaptive visual tracking for robotic systems," *IEEE Trans. Control Syst. Technol.*, vol. 26, no. 6, pp. 2232–2241, Apr. 2018.
- [5] Y. Li, H. Wang, Y. Xie, C. C. Cheah, and W. Ren, "Adaptive image-space regulation for robotic systems," *IEEE Trans. Control Syst. Technol.*, vol. 29, no. 2, pp. 850–857, Dec. 2021.
- [6] Y. Zhang and C. Hua, "A new adaptive visual tracking scheme for robotic system without image-space velocity information," *IEEE Trans. Syst., Man, Cybern.: Syst.*, vol. 52, no. 8, pp. 5249–5258, Aug. 2022.
- [7] G. Hu, S. Gupta, N. Fitz-Coy, and W. E. Dixon, "Lyapunov-based visual servo tracking control via a quaternion formulation," in *Proc. IEEE Conf. Decis. Control*, San Diego, CA, USA, 2006, pp. 3861–3866.
- [8] J. Chen, V. K. Chitrakaran, and D. M. Dawson, "Range identification of features on an object using a single camera," *Automatica*, vol. 47, no. 1, pp. 201–206, Jan. 2011.
- [9] S. Mehta, V. Jayaraman, T. Burks, and W. Dixon, "Teach by zooming: A unified approach to visual servo control," *Mechatronics*, vol. 22, no. 4, pp. 436–443, Jun. 2012.
- [10] Y. Pan and H. Yu, "Composite learning robot control with guaranteed parameter convergence," *Automatica*, vol. 89, pp. 398–406, Mar. 2018.
- [11] B. Lai, Z. Li, W. Li, C. Yang, and Y. Pan, "Homography-based visual servoing of eye-in-hand robots with exact depth estimation," *IEEE Trans. Ind. Electron.*, vol. 71, no. 4, pp. 3832–3841, Apr. 2024.
- [12] Z. Li, B. Lai, and Y. Pan, "Image-based composite learning robot visual servoing with an uncalibrated eye-to-hand camera," *IEEE/ASME Trans. Mechatronics*, vol. 29, no. 4, pp. 2499–2509, Aug. 2024.
- [13] Z. Li, W. Li, and Y. Pan, "Composite learning image-based visual servoing of redundant robots with nullspace compliance," *IEEE Control Syst. Lett.*, vol. 8, pp. 315–320, Jan. 2024.
- [14] Z. Li, B. Lai, H. Wang, and Y. Pan, "Homography-based visual servoing of robot pose under an uncalibrated eye-to-hand camera," *IEEE/ASME Trans. Mechatronics*, vol. 29, no. 3, pp. 1891–1902, Jun. 2024.
- [15] R. Hartley and A. Zisserman, *Multiple View Geometry in Computer Vision*, 2nd ed. New York, NJ, USA: Cambridge Univ. Press, 2004.
- [16] J. J. Craig, *Introduction to Robotics: Mechanics and Control*, 3rd ed. Upper Saddle River, NJ, USA: Prentice-Hall, 2005.
- [17] X. Liang, H. Wang, Y.-H. Liu, W. Chen, and J. Zhao, "A unified design method for adaptive visual tracking control of robots with eye-in-hand/fix camera configuration," *Automatica*, vol. 59, pp. 97–105, Sep. 2015.
- [18] Y.-H. Liu, H. Wang, C. Wang, and K. K. Lam, "Uncalibrated visual servoing of robots using a depth-independent interaction matrix," *IEEE Trans. Robot.*, vol. 22, no. 4, pp. 804–817, Aug. 2006.
- [19] H. Wang, "Passivity-based adaptive control for visually servoed robotic systems," in *Proc. Austral. Control Conf.*, Canberra, Australia, 2014, pp. 152–157.
- [20] J. Wang and E. Olson, "AprilTag 2: Efficient and robust fiducial detection," in *Proc. IEEE/RSJ Int. Conf. Intell. Robot. Syst.*, Daejeon, South Korea, 2016, pp. 4193–4198.
- [21] E. Marchand, F. Spindler, and F. Chaumette, "ViSP for visual servoing: A generic software platform with a wide class of robot control skills," *IEEE Robot. Autom. Mag.*, vol. 12, no. 4, pp. 40–52, Dec. 2005.