

3D Feature-Based Sampled-Data Visual Tracking^{*}

Marco Costanzo^{*} Giuseppe De Maria^{*} Ciro Natale^{*}
Antonio Russo^{*}

^{*} *Dipartimento di Ingegneria,
Università degli Studi della Campania “Luigi Vanvitelli”,
Aversa, 81031 ITALY
(e-mail:*

{marco.costanzo,giuseppe.demaria,ciro.natale,antonio.russo1}@unicampania.it}

Abstract: 3D feature-based Visual Servoing (VS) on the one hand shows attractive peculiarities, on the other hand it suffers from drawbacks related to the existence of local minima, which may affect the convergence character of the VS control loop. Furthermore, the performance of the visual tracking module may constitute a bottleneck enforcing severe constraints on the workspace and visual task execution speed. In this paper we introduce a novel sampled-data model of the 3D feature-based VS, and, in order to avoid drawbacks due to local minima, we plan the target reference trajectory in the feature space with the aim to constraint the feature error dynamics to remain close to the desired equilibrium point. Then, we propose a novel feature generation based on the homography provided by a template matching algorithm based on the Zero mean Normalized Cross Correlation (ZNCC) and the design of a visual tracking scheme by resorting to the Extended Kalman Filter (EKF) and Lyapunov direct method, which explicitly takes into account the camera velocity limits, while ensuring stability.

Copyright © 2023 The Authors. This is an open access article under the CC BY-NC-ND license (<https://creativecommons.org/licenses/by-nc-nd/4.0/>)

Keywords: Robots manipulators, Perception and sensing, Guidance navigation and control, Tracking, Lyapunov methods.

1. INTRODUCTION

The objective of accurately positioning the robot end effector with respect to a target object in dynamic environments is a challenge in robotic vision control. This is due to the scene changing over time, which may generate temporary partial occlusions, variation of object location and lighting conditions. Classical vision controllers (refer to Corke (2017) for an up-to-date textbook on visual servoing), designed to achieve the above objective are based on features extracted from the image acquired by the camera mounted on the end effector (*eye-in-hand* configuration) and can be grouped in two main categories, namely Position Based Visual Servoing (PBVS) and Image Based Visual Servoing (IBVS). The former uses the extracted features to estimate the camera pose directly in the 3D space with the aim to minimize the mismatch between the current and the target pose, while the latter aims at minimizing the error between the current image features and the target ones directly in the feature space. The availability of cheap and lightweight RGB-D cameras promotes the use of 3D features, obtained by combining 2D and depth information (Cervera et al., 2003; Costanzo et al., 2022), which allows the design of 3D feature-based VS. Feature-based VS schemes (IBVS and 3D feature-based VS), differently from PBVS, do not require the 3D pose reconstruction, they are robust against camera calibration errors and, since the image features are directly controlled,

the feature-based VS controller will likely be able to keep them in the camera field of view (Cervera et al., 2003; Costanzo et al., 2022). The features extracted from the current image are matched with the ones of the target image acquired beforehand, and, then the tracking module provides, frame by frame, the feature vector as input to the visual controller. The tracking module is the core of any vision control scheme and it should be accurate and robust against lighting conditions and object partial occlusions and motion, as well as it should guarantee alignment accuracy of the current image with the target one. Tracking methods can be feature-based, i.e. they track geometrical primitives or object contour, or model-based which use a model of the tracked object. Among the latter methods, we refer to those based on 2D models that are reference template on the target image. The objective is to estimate the warp between the current image and the reference one. The estimation is obtained by a nonlinear optimization process based on a measure of the alignment between the two images. The result of the estimation process can be either linear translation parameters or an homography, which links the current and the reference image. The model-based approaches differ on the measure function used in the estimation process. The standard one is the Sum of Squared Differences (SSD) function (Lucas and Kanade, 1981; Baker and Matthews, 2004; Benhimane and Malis, 2007), but, unfortunately, this approach is very sensitive to lighting conditions and occlusions. More robust methods are those using similarity measures invariant to changes in brightness and contrast, such as Mutual

^{*} This work was partially supported by the European Commission within H2020 REMODEL Project (GA no. 870133).

Information (MI) (Dame and Marchand, 2010) or Zero mean Normalized Cross Correlation (ZNCC) (Di Stefano et al., 2005). Even if MI presents a significant optimum, its convergence basin is very small with respect to ZNCC and this makes it sensitive to the camera velocity. In this paper we use the ZNCC measure to estimate the homography, which links the current image and the reference template designed by the user on the image of a planar target. Then, the homography is used to generate, frame by frame, artificial features on the current image, matched with the ones arbitrarily selected by the user within the reference template, as input to the visual controller.

In order to plan the motion of the camera in terms of displacement, velocity and acceleration and to avoid an anomalous character of the feature trajectories, discussed in Costanzo et al. (2022), due to the existence of local minima, we design a time-varying reference feature trajectory starting from the features selected within the template on the target image and the ones generated by means of the homography estimated in the initial condition. If the visual servoing task starts with a feature error equal to zero, then the feature error dynamics is constrained to stay close to the desired equilibrium point. Arbitrary motion of the target object, (e.g., an object handled by a human during an handover maneuver or a picking operation of an object on a conveyor belt), may affect the performance of the visual tracker. To avoid this drawback, we estimate the object velocity by means of an Extended Kalman Filter (EKF) based on a sampled-data model of the VS introduced in this paper without resorting to the forward Euler approximation. Then, we add it to the camera velocity command to compensate the relative motion, not related to the reference feature, between the camera and the object, thus significantly reducing tracker sensitivity to the object velocity. Finally, we propose a visual controller by resorting to the Lyapunov direct method, which explicitly takes into account the camera velocity limits, while ensuring stability at the same time.

2. 3D FEATURE-BASED VISUAL SERVOING SAMPLED DATA MODEL

The VS consists in the design of a closed-loop control scheme that regulates the current image feature set to the given target of an image acquired beforehand. Since the sampling rate that can be achieved by the visual perception pipeline is much lower than the one of the low-level robot control loop, the sampled-data nature of the visual servoing should be taken into account in order to allow maximum possible performance of the VS task.

Let $\mathbf{s}_i, \mathbf{s}_{f_i}^*, \mathbf{s}_i^* \in \mathbb{R}^3, i = 1, \dots, n$ be the i th 3D-point feature vectors of the current, target and reference image, respectively, and let $\mathbf{s} = [\mathbf{s}_1^\top \dots \mathbf{s}_n^\top]^\top \in \mathbb{R}^{3n}$, $\mathbf{s}_f^* = [\mathbf{s}_{f_1}^{*\top} \dots \mathbf{s}_{f_n}^{*\top}]^\top \in \mathbb{R}^{3n}$ be the vectors of n matched image features. The reference features $\mathbf{s}_i^*(t)$ are designed starting from the initial features $\mathbf{s}_i(0)$ to the target one $\mathbf{s}_{f_i}^*$ to plan the feature trajectory during the visual servoing task.

Denote with $\mathbf{v} = [\mathbf{v}^\top \boldsymbol{\omega}^\top]^\top \in \mathbb{R}^6$ the body velocity screw of the camera, where \mathbf{v} is the linear velocity and $\boldsymbol{\omega}$ is the angular velocity. The relationship between the rate of $\mathbf{s}(t)$ and the velocity screw $\mathbf{v}(t)$ is described as

$$\dot{\mathbf{s}}(t) = \mathbf{L}(\mathbf{s})\mathbf{v}(t), \quad \mathbf{s}(0) = \mathbf{s}_0, \quad (1)$$

where $\mathbf{L}(\mathbf{s}) \in \mathbb{R}^{3n \times 6}$ is the so-called *interaction matrix* which, due to the particular feature set selected, takes the structure

$$\mathbf{L}(\mathbf{s}) = \begin{bmatrix} -\mathbf{I}_3 & \dots & -\mathbf{I}_3 \\ \mathbf{S}^\top(\mathbf{s}_1) & \dots & \mathbf{S}^\top(\mathbf{s}_n) \end{bmatrix}^\top, \quad (2)$$

where \mathbf{I}_3 indicates the 3×3 identity matrix and $\mathbf{S}(\mathbf{s}_i)$ is the skew symmetric operator applied to the vector \mathbf{s}_i , such that $\mathbf{S}(\mathbf{s}_i)\boldsymbol{\omega} = \mathbf{s}_i \times \boldsymbol{\omega}$.

By defining the feature error vector $\mathbf{e} \in \mathbb{R}^{3n}$ as

$$\mathbf{e}(t) = \mathbf{s}(t) - \mathbf{s}^*(t), \quad (3)$$

denoting with $\mathbf{L}_e = \mathbf{L}(\mathbf{e} + \mathbf{s}^*)$, the dynamics of the current feature error can be easily derived from (1) and (3) as

$$\dot{\mathbf{e}}(t) = \mathbf{L}_e \mathbf{v}(t) - \dot{\mathbf{s}}^*(t), \quad \mathbf{e}(0) = \mathbf{e}_0. \quad (4)$$

In the following, the dependence from the time will be dropped to simplify the notation.

Taking into account the structure of $\mathbf{L}(\mathbf{s})$ in (2), the term at the right-hand side of equation (1) can be written as $\mathbf{L}(\mathbf{s})\mathbf{v} = \bar{\mathbf{S}}^\top(\boldsymbol{\omega})\mathbf{s} - \bar{\mathbf{v}}$ and the feature dynamics can be expressed as

$$\dot{\mathbf{s}} = \bar{\mathbf{S}}^\top(\boldsymbol{\omega})\mathbf{s} - \bar{\mathbf{v}}, \quad (5)$$

where $\bar{\mathbf{S}}(\boldsymbol{\omega}) \in \mathbb{R}^{3n \times 3n}$ and $\bar{\mathbf{v}} \in \mathbb{R}^{3n}$ have the following expressions

$$\bar{\mathbf{S}}(\boldsymbol{\omega}) = \mathbf{I}_n \otimes \mathbf{S}(\boldsymbol{\omega}), \quad \bar{\mathbf{v}} = \mathbf{1}_n \otimes \mathbf{v}. \quad (6)$$

where the symbol \otimes indicates the Kronecker product, \mathbf{I}_n is the $n \times n$ identity matrix and $\mathbf{1}_n = [1 \ 1 \ \dots \ 1]^\top \in \mathbb{R}^n$.

Denoting with T the control sampling time, the camera velocity control is kept constant at the value $\mathbf{v}(kT) = \mathbf{v}_k = [\mathbf{v}_k^\top \boldsymbol{\omega}_k^\top]^\top$ in the time interval $[kT, (k+1)T)$. Therefore, $\bar{\mathbf{S}}(\boldsymbol{\omega}_k)$ is updated at each step k and is constant within this time interval. The feature dynamics can be computed as the solution of (5) in $[kT, (k+1)T)$, i.e.,

$$\mathbf{s}(kT+t) = e^{-\bar{\mathbf{S}}(\boldsymbol{\omega}_k)t} \mathbf{s}(kT) - \int_{kT}^{kT+t} e^{-\bar{\mathbf{S}}(\boldsymbol{\omega}_k)(kT+t-\tau)} d\tau \bar{\mathbf{v}}_k \quad (7)$$

for $t \in [0, T)$. The feature dynamics in (7), in the discrete time k , is similar to that of switched systems. Relying on the results in Liberzon and Morse (1999), for $t = T$ and by a suitable change of variable, the feature dynamics in the sampling instants is governed by

$$\mathbf{s}_{k+1} = e^{-\bar{\mathbf{S}}(\boldsymbol{\omega}_k)T} \mathbf{s}_k - \int_0^T e^{-\bar{\mathbf{S}}(\boldsymbol{\omega}_k)\sigma} d\sigma \bar{\mathbf{v}}_k, \quad \forall k \in \mathbb{N}_0, \quad (8)$$

where \mathbf{s}_k stands for $\mathbf{s}(kT)$.

By expanding the exponential matrices, in the previous equation, in power series, it holds that

$$\int_0^T e^{-\bar{\mathbf{S}}(\boldsymbol{\omega}_k)\sigma} d\sigma = T \sum_{i=0}^{\infty} \frac{(\bar{\mathbf{S}}^\top(\boldsymbol{\omega}_k))^i T^i}{(i+1)!}, \quad (9)$$

thus, substituting (9) in (8) and rearranging the terms, equation (8) becomes

$$\mathbf{s}_{k+1} = \mathbf{s}_k + \mathbf{P}(\boldsymbol{\omega}_k)\mathbf{L}(\mathbf{s}_k)\mathbf{v}_k, \quad (10)$$

with $\mathbf{P}(\cdot) : \mathbb{R}^3 \rightarrow \mathbb{R}^{3n \times 3n}$ defined as

$$\mathbf{P}(\boldsymbol{\omega}_k) = \int_0^T e^{-\bar{\mathbf{S}}(\boldsymbol{\omega}_k)\sigma} d\sigma. \quad (11)$$

Equation (10) represents the sampled-data feature dynamic model, which describes exactly the behaviour of the VS system at the sampling instants.

Now, by denoting with $\mathbf{e}_k = \mathbf{e}(kT)$, then it is straightforward to derive the feature error dynamics as

$$\mathbf{e}_{k+1} = \mathbf{e}_k + \mathbf{P}(\boldsymbol{\omega}_k) \mathbf{L}_{\mathbf{e}_k} \mathbf{v}_k - \Delta \mathbf{s}_k^*, \quad (12)$$

where $\mathbf{L}_{\mathbf{e}_k} = \mathbf{L}(\mathbf{e}_k + \mathbf{s}^*)$ and $\mathbf{v}_k = [\mathbf{v}_k^\top \boldsymbol{\omega}_k^\top]^\top \in \mathbb{R}^6$ and $\Delta \mathbf{s}_k^* = \mathbf{s}_{k+1}^* - \mathbf{s}_k^*$ is the forward finite difference of \mathbf{s}_k^* .

3. TIME-VARYING REFERENCE FEATURE GENERATION

The time-varying reference $\mathbf{s}^*(t)$ is designed to keep the feature error close to the origin during the whole task as to avoid the undesired behaviours discussed in (Costanzo et al., 2022). Such a reference can be generated by interpolating each feature between the initial value $\mathbf{s}(0)$ and the target one \mathbf{s}_f^* in a given time t_f . The interpolation is obtained by applying the time-varying rigid transformation $\mathbf{T}^*(t) = \mathbf{T}(\theta^*(t), \mathbf{r}^*(t), \mathbf{p}^*(t)) \in \mathbb{R}^{4 \times 4}$ representing a rotation θ^* about a rotation axis \mathbf{r}^* , and a translation \mathbf{p}^* , to each feature. Thus, the reference features are obtained as

$$\tilde{\mathbf{s}}_i^*(t) = \mathbf{T}^*(t) \tilde{\mathbf{s}}_i(0), \quad i = 1, \dots, n, \quad (13)$$

where $\tilde{\mathbf{s}}_i^*$ and $\tilde{\mathbf{s}}_i$ indicate the vector of homogeneous coordinates of the feature \mathbf{s}_i^* and \mathbf{s}_i , respectively.

Then, assuming $\mathbf{r}^*(t) = \mathbf{r}_f^*$ constant, the designed homogeneous transformation $\mathbf{T}(\theta^*, \mathbf{r}^*, \mathbf{p}^*)$ should be such that

$$\mathbf{T}(\theta^*(0), \mathbf{r}_f^*, \mathbf{p}^*(0)) = \mathbf{T}^*(0) = \mathbf{I}_4, \quad (14)$$

$$\mathbf{T}(\theta^*(t_f), \mathbf{r}_f^*, \mathbf{p}^*(t_f)) \tilde{\mathbf{s}}_i(0) = \mathbf{T}^* \tilde{\mathbf{s}}_i(0) = \tilde{\mathbf{s}}_i^*. \quad (15)$$

In the following we design θ^* and \mathbf{p}^* so that they evolve in time as a fifth-order polynomial between $\theta^*(0) = 0$, $\mathbf{p}^*(0) = 0$ and $\theta^*(t_f)$, $\mathbf{p}^*(t_f)$, respectively.

Recalling (15), the homogeneous transformation matrix $\mathbf{T}^*(t_f)$ can be found by solving the following problem

$$\min_{\mathbf{T}} \frac{1}{2} \sum_{i=1}^n \|\mathbf{T} \tilde{\mathbf{s}}_i(0) - \tilde{\mathbf{s}}_i^*\|^2. \quad (16)$$

To ensure that \mathbf{T} is an homogeneous transformation matrix, we parametrize it with a position vector and a unit quaternion.

4. HOMOGRAPHY-BASED 3D FEATURE GENERATION

Measurement of the 3D features can be accomplished by tracking 2D features on the camera image in pixel coordinates and then transforming them in the 3D space via the RGB-D camera intrinsic parameters. We assume that the object to be tracked has a planar textured face on which the features are located (the generalization to the non-planar case is straightforward (Benhimane and Malis, 2007)). Let $[u_{ik} \ v_{ik}]^\top$, $[u_{fi}^* \ v_{fi}^*]^\top$ be the projection of the current 3D feature \mathbf{s}_{ik} and target \mathbf{s}_{fi}^* on the image plane in pixel coordinates, respectively. Since the 3D features are located on a planar face, the relation between the features in the image plane is given by the homography matrix $\mathbf{H}_k \in \mathbb{R}^{3 \times 3}$ (Benhimane and Malis, 2007), i.e.,

$$[\mu_{ik} u_{ik} \ \mu_{ik} v_{ik} \ \mu_{ik}]^\top = \mathbf{H}_k [u_{fi}^* \ v_{fi}^* \ 1]^\top \quad (17)$$

where μ_{ik} is an auxiliary variable.

In our approach, we estimate the homography \mathbf{H}_k via a template-tracking algorithm based on ZNCC (Di Stefano et al., 2005). The algorithm needs to be initialized by selecting the template on the target image and by determining an estimation of the initial homography \mathbf{H}_0 . To estimate \mathbf{H}_0 , we do as follows. Given the initial and target image, couples of 2D features are matched by means of the matching algorithm provided by the ViSP library (Marchand et al., 2005) (vpKeyPoint class). The matched features are used to solve (17) for \mathbf{H}_0 by using the Ransac algorithm (vpHomography class). Then, the adopted tracker (vpTemplateTrackerZNCCInverseCompositional class) solves, at each time step, the following optimization problem to find the homography parameters (Dame and Marchand, 2010)

$$\hat{\mathbf{h}}_k = \arg \max_{\mathbf{h}} f(\mathcal{I}^*, w(\mathcal{I}_k, \mathbf{h})) \quad (18)$$

where $\mathbf{h} \in \mathbb{R}^8$ are the free parameters of the homography matrix \mathbf{H}_k , \mathcal{I}^* and \mathcal{I}_k are the target image template and the current one, respectively, $f(\cdot, \cdot)$ is the similarity function that checks the similarity of two images via the ZNCC algorithm, and $w(\cdot, \cdot)$ is the warping function that transforms the template by applying the homography to each pixel as in (17).

For the control purpose, we select the target features $[u_{fi}^* \ v_{fi}^*]^\top$ arbitrarily within the target template, then we generate the matched current features $[u_{ik} \ v_{ik}]^\top$ via the homography (17).

5. OBJECT VELOCITY ESTIMATION

The object motion generates an additional velocity \mathbf{v}_d with respect to the camera frame, i.e., $\mathbf{v} = \mathbf{v}_c + \mathbf{v}_d$, where \mathbf{v}_c is the actual camera velocity command. By denoting with $\mathbf{v}_c = [\mathbf{v}_c^\top \boldsymbol{\omega}_c^\top]^\top$ and $\mathbf{v}_d = [\boldsymbol{\eta}^\top \boldsymbol{\zeta}^\top]^\top$, where $\boldsymbol{\eta} \in \mathbb{R}^3$ and $\boldsymbol{\zeta} \in \mathbb{R}^3$ are the linear and angular components of the object velocity, respectively, \mathbf{v}_d can be estimated by resorting to an EKF tailored on the following equations where a constant object velocity model is assumed,

$$\mathbf{v}_{dk+1} = \mathbf{v}_{dk} + \boldsymbol{\nu}_k, \quad (19)$$

$$\begin{aligned} \mathbf{s}_{k+1} &= \mathbf{s}_k + \mathbf{P}(\boldsymbol{\omega}_{ck} + \boldsymbol{\zeta}_k) \mathbf{L}(\mathbf{s}_k) (\mathbf{v}_{ck} + \mathbf{v}_{dk}) + \boldsymbol{\sigma}_k, \\ &= \mathbf{f}(\mathbf{s}_k, \boldsymbol{\eta}_k, \boldsymbol{\zeta}_k) + \boldsymbol{\sigma}_k \end{aligned} \quad (20)$$

$$\mathbf{y}_k = \mathbf{s}_k + \boldsymbol{\chi}_k, \quad (21)$$

where $\boldsymbol{\nu}_k$, $\boldsymbol{\sigma}_k$, $\boldsymbol{\chi}_k$ are Gaussian process noise with covariance matrices $\boldsymbol{\Gamma} \in \mathbb{R}^{6 \times 6}$, $\boldsymbol{\Sigma} \in \mathbb{R}^{3n \times 3n}$ and $\mathbf{X} \in \mathbb{R}^{3n \times 3n}$, respectively, and \mathbf{y}_k is the system measurable output. Computation of the Jacobian matrix \mathbf{J} of system (19)-(20) is not trivial and it is obtained as

$$\mathbf{J} = \begin{bmatrix} \mathbf{I}_3 & \mathbf{0}_3 & \mathbf{0}_3 \\ \mathbf{0}_3 & \mathbf{I}_3 & \mathbf{0}_3 \\ \frac{\partial \mathbf{f}}{\partial \boldsymbol{\eta}} & \begin{bmatrix} \frac{\partial \mathbf{f}}{\partial \zeta_1} & \frac{\partial \mathbf{f}}{\partial \zeta_2} & \frac{\partial \mathbf{f}}{\partial \zeta_3} \end{bmatrix} & \frac{\partial \mathbf{f}}{\partial \mathbf{s}} \end{bmatrix}, \quad (22)$$

where

$$\frac{\partial \mathbf{f}}{\partial \boldsymbol{\eta}} = -\mathbf{P}(\boldsymbol{\omega}_c + \boldsymbol{\zeta}) (\mathbf{1}_n \otimes \mathbf{I}_3), \quad \frac{\partial \mathbf{f}}{\partial \mathbf{s}} = e^{-\bar{\mathbf{s}}(\boldsymbol{\omega}_c + \boldsymbol{\zeta})},$$

$$\frac{\partial \mathbf{f}}{\partial \zeta_i} = \frac{\partial \mathbf{P}(\boldsymbol{\omega}_c + \boldsymbol{\zeta})}{\partial \zeta_i} \mathbf{L}(\mathbf{s}) \begin{bmatrix} \mathbf{v}_c + \boldsymbol{\eta} \\ \boldsymbol{\omega}_c + \boldsymbol{\zeta} \end{bmatrix} + \mathbf{P}(\boldsymbol{\omega}_c + \boldsymbol{\zeta}) \mathbf{L}(\mathbf{s}) \begin{bmatrix} \mathbf{0}_3 \\ \mathbf{l}_i \end{bmatrix}$$

where \mathbf{l}_i is the i th canonical orthonormal basis versor.

In order to compute the terms $\partial \mathbf{P} / \partial \zeta_i$, recalling that $\boldsymbol{\omega} = \boldsymbol{\omega}_c + \boldsymbol{\zeta}$, by virtue of the Euler-Rodrigues rotation formula

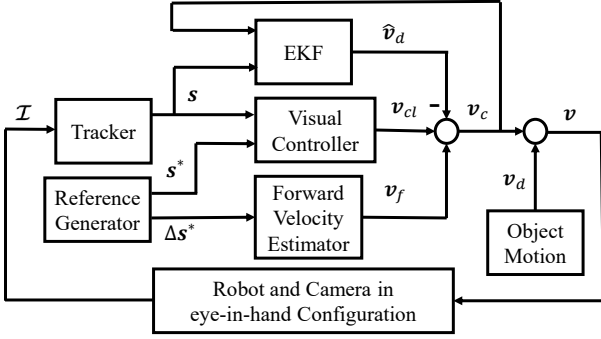


Fig. 1. Visual control system block scheme.

(Dai (2015)), matrix $\mathbf{P}(\boldsymbol{\omega}_k)$ in the system dynamics (20) can be written as

$$\mathbf{P}(\boldsymbol{\omega}_k) = T \left[\mathbf{I}_{3n} + \frac{1 - \cos(\|\boldsymbol{\omega}_k\|T)}{\|\boldsymbol{\omega}_k\|T} \bar{\mathbf{S}}^\top \left(\frac{\boldsymbol{\omega}_k}{\|\boldsymbol{\omega}_k\|} \right) + \frac{\|\boldsymbol{\omega}_k\|T - \sin(\|\boldsymbol{\omega}_k\|T)}{\|\boldsymbol{\omega}_k\|T} \bar{\mathbf{S}}^2 \left(\frac{\boldsymbol{\omega}_k}{\|\boldsymbol{\omega}_k\|} \right) \right]. \quad (23)$$

From (23) it holds that

$$\frac{\partial \mathbf{P}}{\partial \zeta_i} = T \left[\phi'_1(\|\boldsymbol{\omega}\|) \omega_i \bar{\mathbf{S}}^\top(\boldsymbol{\omega}) + \phi_1(\|\boldsymbol{\omega}\|) \bar{\mathbf{S}}^\top(\mathbf{l}_i) + \phi'_2(\|\boldsymbol{\omega}\|) \omega_i \bar{\mathbf{S}}^2(\boldsymbol{\omega}) + \phi_2(\|\boldsymbol{\omega}\|) (\bar{\mathbf{S}}(\mathbf{l}_i) \bar{\mathbf{S}}(\boldsymbol{\omega}) + \bar{\mathbf{S}}(\boldsymbol{\omega}) \bar{\mathbf{S}}(\mathbf{l}_i)) \right]$$

where

$$\phi_1(x) = \sum_{i=1}^{\infty} \frac{(-1)^i x^{2i-2} T^{2i-1}}{(2i)!}, \quad \phi_2(x) = \sum_{i=1}^{\infty} \frac{(-1)^i x^{2i-2} T^{2i}}{(2i+1)!},$$

with $\phi'_j(x) = \frac{1}{x} \frac{d\phi_j}{dx}$, $j = 1, 2$.

The effectiveness of the estimation by the EKF will be evaluated in Section 7.

6. VISUAL TRACKING

Figure 1 shows the control system architecture. The camera velocity command \mathbf{v}_c is obtained as sum of three terms

$$\mathbf{v}_c = -\hat{\mathbf{v}}_d + \mathbf{v}_f + \mathbf{v}_{cl}. \quad (24)$$

The first term $\hat{\mathbf{v}}_d$ is the estimation of the object velocity by the EKF described in Section 5, its purpose is to counteract the object velocity \mathbf{v}_d in the block scheme; \mathbf{v}_f represents the forward control action which is designed to balance the term $\Delta \mathbf{s}_k^*$ in (12) and can be computed by solving the following implicit equation

$$\mathbf{v}_{fk} = \begin{bmatrix} \mathbf{v}_{fk} \\ \boldsymbol{\omega}_{fk} \end{bmatrix} : \mathbf{P}(\boldsymbol{\omega}_{fk}) \mathbf{L}_{e_k} \mathbf{v}_{fk} = \Delta \mathbf{s}_k^*. \quad (25)$$

This equation can be solved by reformulating the problem as a minimization one as follows

$$\min_{\mathbf{v}_{fk}} \|\mathbf{P}(\boldsymbol{\omega}_{fk}) \mathbf{L}_{e_k} \mathbf{v}_{fk} - \Delta \mathbf{s}_k^*\|^2. \quad (26)$$

The last term in (24) is the closed loop control action \mathbf{v}_{cl} and it is designed to ensure the stability of the VS scheme. Such a component will be derived on system (12) assuming that the estimated relative velocity cancels the actual one ($\hat{\mathbf{v}}_d = \mathbf{v}_d$). Thus (12), in view of (25), reduces to

$$\mathbf{e}_{k+1} = \mathbf{e}_k + \mathbf{P}(\boldsymbol{\omega}_{clk}) \mathbf{L}_{e_k} \mathbf{v}_{clk}. \quad (27)$$

With the aim to design a velocity control law that aims at guaranteeing a decoupled and exponential decrease of the error in the case of pure translational motion (Corke, 2017; Costanzo et al., 2022), we assume for the system in (27)

a reference model $\mathbf{e}_{k+1} = (1 - \lambda_k T) \mathbf{e}_k$ with $0 < \lambda_k T < 1$, which assures in the case of pure translational motion a decoupled and convergent dynamics without hidden oscillations. The equivalent control \mathbf{v}_{clk} , which assures a perfect model reference following in pure translation, is $\mathbf{v}_{clk} = -\lambda_k \mathbf{L}_{e_k}^\dagger \mathbf{e}_k$ where λ_k is the control gain. Thus, the closed-loop dynamics becomes

$$\mathbf{e}_{k+1} = (\mathbf{I}_{3n} - \lambda_k \mathbf{P}_k \mathbf{L}_{e_k} \mathbf{L}_{e_k}^\dagger) \mathbf{e}_k. \quad (28)$$

The control law proposed in the following proposition requires the definition of two scalar functions, i.e., $\pi(\cdot, \cdot)$, the *cost function*, and $l(\cdot)$, the *landing function*. The former is denoted as

$$\pi(\mathbf{e}_k, \lambda_k T) = -\lambda_k T (2l(\mathbf{e}_k) - \alpha_k \lambda_k T \|\boldsymbol{\omega}_k^1\| - \lambda_k T), \quad (29)$$

where $\boldsymbol{\omega}_k^1 = -[\mathbf{O}_3 \ \mathbf{I}_3] \mathbf{L}_{e_k}^\dagger \mathbf{e}_k$, $\alpha_k = \|\mathbf{e}_k\| / \|\mathbf{e}_{Ik}\|$ and $\mathbf{e}_{Ik} = \mathbf{L}_{e_k} \mathbf{L}_{e_k}^\dagger \mathbf{e}_k$ is the projection of \mathbf{e}_k in $\text{Im}(\mathbf{L}_e)$. The function $\pi(\mathbf{e}_k, \lambda_k T)$, as a function of $\lambda_k T$, is a concave up parabola passing through the origin and the point $\lambda_k T = 2l(\mathbf{e}_k) / (1 + \alpha_k \|\boldsymbol{\omega}_k^1\|)$.

The *landing function* $l(\mathbf{e}_k)$, which allows a smooth approach to the reference trajectory, is defined as follows. Denote with $e_{Mk} = \|\mathbf{e}_k\| / \sqrt{n}$ the Root Mean Square Error (RMSE), let e_H and e_L be two scalar values such that $e_L < e_H$ and $\beta \in (0, 1]$, then

$$l(\mathbf{e}_k) = \begin{cases} \beta & \text{if } e_{Mk} \leq e_L \\ 1 & \text{if } e_{Mk} \geq e_H \\ \beta + (1 - \beta) q\left(\frac{e_{Mk} - e_L}{e_H - e_L}\right) & \text{otherwise} \end{cases}, \quad (30)$$

where $q(\cdot)$ is a fifth-order polynomial whose coefficients are determined such that the first and second order derivatives are zero at the extremes of the interval $[0, 1]$ and $q(0) = 0$, $q(1) = 1$. Note that $0 < l(\mathbf{e}_k) \leq 1$ and $\beta = 1$ implies $l(\mathbf{e}_k) = 1$.

Now denote with ω_k^M and v_k^M the maximum rotational and linear camera velocities, respectively, depending on the selected joint velocity limits and the configuration of the robot performing the task. Moreover, let $\mathcal{D} \subset \mathbb{R}^{3n}$ be an open neighborhood of the equilibrium point $\mathbf{e} = \mathbf{0}$, which does not contain any other point satisfying the condition $\mathbf{L}_e^\top \mathbf{e} = \mathbf{0}$ (see (Costanzo et al., 2022) for details). The following proposition holds.

Proposition 1. If the control gain λ_k is determined for each sampling interval by solving the following problem

$$\min_{\lambda_k T} \pi(\mathbf{e}_k, \lambda_k T) \quad (31a)$$

$$\text{s.t. } 0 < \lambda_k T < \frac{2l(\mathbf{e}_k)}{1 + \alpha_k \|\boldsymbol{\omega}_k^1\|} \quad (31b)$$

$$\|\boldsymbol{\omega}_k\| \leq \omega^M, \quad \|\mathbf{v}_k\| \leq v^M, \quad (31c)$$

then the equilibrium point $\mathbf{e} = \mathbf{0}$ is asymptotically stable in \mathcal{D} and the feature error dynamics does not exhibit hidden oscillations over the continuous time t .

Proof. By selecting as Lyapunov function candidate $V(\mathbf{e}_k) = \mathbf{e}_k^\top \mathbf{e}_k$, its first order variation $\Delta V_k = \mathbf{e}_{k+1}^\top \mathbf{e}_{k+1} - \mathbf{e}_k^\top \mathbf{e}_k$ along the error dynamics (28) can be written, by denoting $\mathbf{B}_k = \mathbf{P}_k - T \mathbf{I}_{3n}$ with $\mathbf{P}_k = \mathbf{P}(\boldsymbol{\omega}_k)$, as

$$\Delta V_k = -2\lambda_k T \mathbf{e}_k^\top \mathbf{L}_{e_k} \mathbf{L}_{e_k}^\dagger \mathbf{e}_k - 2\lambda_k \mathbf{e}_k^\top \mathbf{B}_k \mathbf{L}_{e_k} \mathbf{L}_{e_k}^\dagger \mathbf{e}_k + \lambda_k^2 \mathbf{e}_k^\top \mathbf{L}_{e_k} \mathbf{L}_{e_k}^\dagger \mathbf{P}_k^\top \mathbf{P}_k \mathbf{L}_{e_k} \mathbf{L}_{e_k}^\dagger \mathbf{e}_k. \quad (32)$$

By neglecting high-order terms in the series expansions of functions in (23), the matrix \mathbf{P}_k can be approximated as

$$\mathbf{P}(\boldsymbol{\omega}_k) \approx T\mathbf{I}_{3n} + \frac{1}{2}T^2\bar{\mathbf{S}}^\top(\boldsymbol{\omega}_k) \quad (33)$$

thus it results $\|\mathbf{P}_k^\top \mathbf{P}_k\| \leq T^2$ and $\|\mathbf{B}_k\| \leq \|\boldsymbol{\omega}_k\|T^2/2$. Since α_k is finite in \mathcal{D} and it is lowerbounded by 1 (Costanzo et al., 2022), it is easy to show that ΔV_k can be overbounded as

$$\Delta V_k \leq -\lambda_k T (2 - \lambda_k T \alpha_k \|\boldsymbol{\omega}_k^1\| - \lambda_k T) \|e_{Ik}\|^2 \quad (34)$$

implying, in view of (29) and the fact that $l(e_k) \in (0, 1]$

$$\Delta V_k \leq \pi(e_k, \lambda_k T) \|e_{Ik}\|^2. \quad (35)$$

Then asymptotic stability in \mathcal{D} is achieved for the desired equilibrium point if, for each k , λ_k satisfies (31b). Using (32) and (33) it is easy to show that $\Delta V_k / \|e_{Ik}^2\|$ as function of $\lambda_k T$ is a concave upward parabola since the last term in (32) is positive definite in \mathcal{D} . For each k , the function $\pi(e_k, \lambda_k T)$ has a minimum for $\lambda_k T = \lambda_k^* T$, which is lower than the one for which the function $\Delta V_k / \|e_{Ik}\|^2$ exhibits its minimum value. Then, for $\lambda_k T \leq \lambda_k^* T$, it is very easy to prove that $\|e(kT + t)\|$ for $t \in [0, T)$ is monotonically decreasing and the feature error dynamics does not exhibit hidden oscillations. If $\lambda_k = \lambda_k^*$ is such that the constraint (31c) is not satisfied the minimization algorithm selects $\lambda_k < \lambda_k^*$ as the camera velocity depends linearly on λ_k . \square

Remark 1. If the EKF estimation is not perfect, by denoting with $\tilde{\mathbf{v}}_{dk} = [\tilde{\eta}_k^\top \tilde{\zeta}_k^\top]^\top = \mathbf{v}_{dk} - \hat{\mathbf{v}}_{dk}$ and folding it into the feature dynamics (10), the error dynamics in (27) can be rewritten as

$$\mathbf{e}_{k+1} = (\mathbf{I}_{3n} - \lambda_k \mathbf{P}_k \mathbf{L}_{e_k} \mathbf{L}_{e_k}^\dagger) \mathbf{e}_k + \boldsymbol{\delta}_k, \quad (36)$$

where $\boldsymbol{\delta}_k = \boldsymbol{\delta}(e_k, \tilde{\mathbf{v}}_{dk})$ is the cumulative function of all the terms depending on $\tilde{\mathbf{v}}_{dk}$, that is limited if the EKF converges. Due to the estimation error, the feature error does not converge to zero but it remains bounded in a neighborhood of $\mathbf{e} = \mathbf{0}$. To prove the existence of such a bound, we derive ΔV_k by resorting to (33), then an upperbound of ΔV_k can be obtained as follows

$$\Delta V_k \leq -\lambda_k T (2 - \lambda_k T \alpha_k \|\boldsymbol{\omega}_k^1\| - \lambda_k T) \|e_{Ik}\|^2 + 2(1 + 2\lambda_k T) \|\boldsymbol{\gamma}_k\| \|e_{Ik}\| + \|\boldsymbol{\gamma}_k\|^2, \quad (37)$$

where terms of order higher than T^2 have been neglected without affecting the sign of the right-hand side, and $\boldsymbol{\gamma}_k = \boldsymbol{\gamma}(\tilde{\mathbf{v}}_{dk})$ is a function depending linearly on $\tilde{\mathbf{v}}_{dk}$ only (space limitations do not allow us to provide all the details of ΔV_k derivation). This upperbound of ΔV_k , as function of $\|e_{Ik}\|$ with λ_k satisfying constraint (31b), is a concave downward parabola which has only one non-negative root e_{Ik}^* . The greater $\|\boldsymbol{\gamma}_k\|$, the greater e_{Ik}^* . Thus, for $\|e_{Ik}\| > e_{Ik}^*$, ΔV_k is negative definite. The better the EKF estimation, the lower the bound e_{Ik}^* .

7. EXPERIMENTS

The proposed controller has been experimentally tested by using a Kuka LBR iiwa 7 robot equipped with an Intel RealSense D435i RGB-D camera in an *eye-in-hand* configuration. The target image is the one of a plastic bottle acquired beforehand (Figure 2-right). The target features have been selected by the user on the bottle surface (right image), while the current features (left

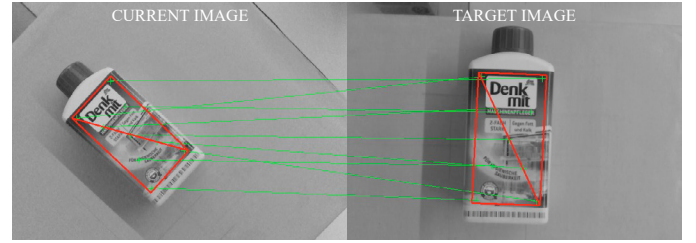


Fig. 2. Target features (right) and matched current ones (left). The area inside the red triangles represents the tracked template.

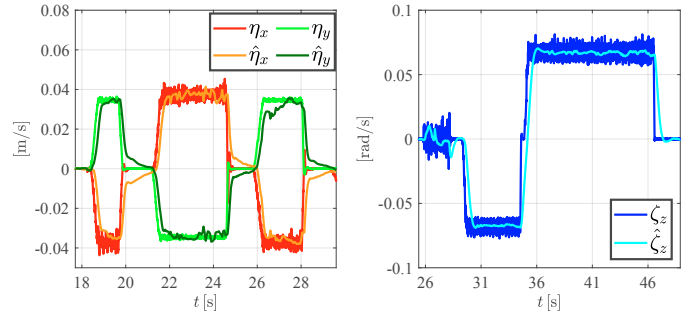


Fig. 3. Exogenous velocity estimated by the EKF: translational (left) and rotational (right).

image) are obtained by applying the homography (17) to the target ones.

The first experiment is devoted to test the accuracy of the object velocity estimation by the EKF. To this aim, we open the loop at the output of the first sum node in Fig. 1 and then move the robot using the teach pendant while keeping the object fixed. This way, the velocity \mathbf{v}_d is the measured robot Cartesian velocity. The noise covariance matrices are selected as $\boldsymbol{\Gamma} = \text{diag}(5\mathbf{I}_3, 20\mathbf{I}_3) \times 10^{-7}$, $\boldsymbol{\Sigma} = 10^{-9}\mathbf{I}_{3n}$ and $\mathbf{X} = 10^{-7}\mathbf{I}_{3n}$. The results are shown in Fig. 3, the EKF is able to accurately estimate the camera/object relative velocity, in particular in the time intervals where the relative velocity is constant, the estimated one is a filtered version of the actual velocity measured as robot Cartesian velocity.

In the second experiment, the camera starts from the configuration depicted in Fig. 2. The desired feature trajectory \mathbf{s}^* is generated with the procedure explained in Section 3 with $t_f = 6$ s. The forward action velocity \mathbf{v}_f is computed by solving problem (26) and the closed loop control \mathbf{v}_{cl} is obtained as described in Section 6. When the approach to the target is completed (i.e., when $\mathbf{s}^* = \mathbf{s}_f^*$) a human operator perturbs the object pose by applying both translational and rotational velocities. The experiment is repeated two times, in the first time the EKF estimation $\hat{\mathbf{v}}_d$ is not used in the control action (24), while in the second time the estimated velocity is used to counteract the relative velocity \mathbf{v}_d . Figure 4 shows a comparison between the two task executions. The top plot shows the root mean square errors. The robot moves from $\mathbf{s}_0 = \mathbf{s}_0^*$ by following the reference trajectory up to $t = 6$ s, then the object is moved by the human operator. In both executions the error is bounded during the object motion, but, when the EKF estimation is not used (blue line), the bound is 2.5 times higher. The bottom plot shows the control gain

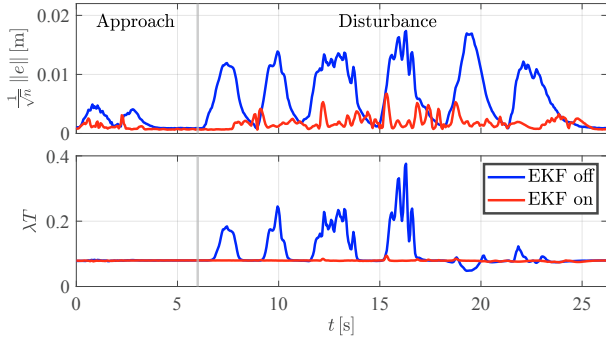


Fig. 4. Second experiment. RMSE (top); control gain λT (bottom).

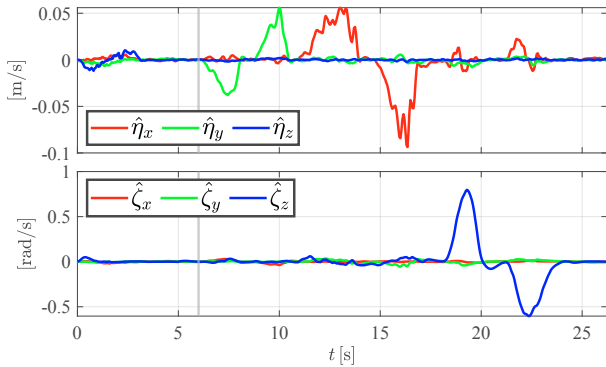


Fig. 5. Second experiment: Estimated velocity in the first task execution.

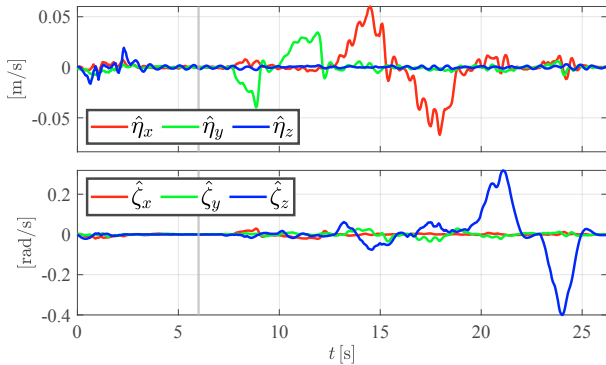


Fig. 6. Second experiment: Estimated velocity in the second task execution.

λ , when the EKF estimation is used (red line) the gain of the closed loop component \mathbf{v}_{cl} is very low since the error is low being $\hat{\mathbf{v}}_d$ an accurate estimation of \mathbf{v}_d . On the other hand, when the EKF estimation is not used (blue line) the control gain increases during the object motion and the camera is controlled exclusively by the closed loop component \mathbf{v}_{cl} and the forward action \mathbf{v}_f . Figure 5 and Fig. 6 show the relative velocity estimated by the EKF during the two task executions. Note that during the approach phase in both executions, the EKF estimates a small relative velocity even if the object is fixed. This is due to the fact that the robot is not a perfect velocity generator and its tracking error can be seen as a disturbance equivalent to a relative object/camera velocity. Also in this phase the error is lower when the EKF estimation is used in the control action.

The video of the experiments is available at <https://youtu.be/beNiXAdrnu>.

8. CONCLUSION

In this paper we have proposed a 3D feature-based sampled-data visual control scheme which takes into account the sampling rate of the visual perception pipeline. To avoid anomalous behaviours due to local minima, we have suitably designed a reference feature trajectory to the control system in order to constraint the feature dynamics close to the desired equilibrium point. To improve the performance of the tracking module in terms of robustness against lighting conditions and partial occlusions, we adopted a ZNCC-based template matching algorithm. The current features are generated by the homography, provided by the visual tracker, applied to the features selected on the target image by the user. As the target object moves, e.g., during an handover maneuver, and this can affect the performance of the visual tracking module, the object motion has been estimated by means of an EKF. Such an estimation has been added to the camera velocity command so that the camera/object relative motion is determined by the reference feature trajectory only. Stability of the closed loop system is ensured by the visual controller, designed by resorting to the Lyapunov direct method, and explicitly taking into account the robot joint velocity limits.

REFERENCES

- Baker, S. and Matthews, I. (2004). Lucas-kanade 20 years on: A unifying framework. *International Journal of Computer Vision*, 56(3), 221–255.
- Benhimane, S. and Malis, E. (2007). Homography-based 2d visual tracking and servoing. *The International Journal of Robotics Research*, 26(7), 661–676.
- Cervera, E., del Pobil, A.P., Berry, F., and Martinet, P. (2003). Improving image-based visual servoing with three-dimensional features. *The International Journal of Robotics Research*, 22(10-11), 821–839.
- Corke, P. (2017). *Robotics, Vision and Control*. Springer Nature, Cham, CH.
- Costanzo, M., De Maria, G., Natale, C., and Russo, A. (2022). Stability and convergence analysis of 3d feature-based visual servoing. *IEEE Robotics and Automation Letters*, 7(4), 12022–12029.
- Dai, J.S. (2015). Euler–Rodrigues formula variations, quaternion conjugation and intrinsic connections. *Mechanism and Machine Theory*, 92, 144–152.
- Dame, A. and Marchand, E. (2010). Accurate real-time tracking using mutual information. In *2010 IEEE Int. Symp. on Mixed and Augmented Reality*, 47–56.
- Di Stefano, L., Mattoccia, S., and Tombari, F. (2005). Zncc-based template matching using bounded partial correlation. *Pattern Recognition Letters*, 26(14), 2129–2134.
- Liberzon, D. and Morse, A. (1999). Basic problems in stability and design of switched systems. *IEEE Control Systems Magazine*, 19(5), 59–70.
- Lucas, B.D. and Kanade, T. (1981). An iterative image registration technique with an application to stereo vision. In *Proc. Int. Joint Conf. on Artificial Intelligence*, 674–679. San Francisco, CA, USA.
- Marchand, E., Spindler, F., and Chaumette, F. (2005). ViSP for visual servoing: a generic software platform with a wide class of robot control skills. *IEEE Robotics Automation Magazine*, 12(4), 40–52.