



# Vision-based Manipulation of Deformable and Rigid Objects Using Subspace Projections of 2D Contours

Jihong Zhu, David Navarro-Alarcon, Robin Passama, Andrea Cherubini

## ► To cite this version:

Jihong Zhu, David Navarro-Alarcon, Robin Passama, Andrea Cherubini. Vision-based Manipulation of Deformable and Rigid Objects Using Subspace Projections of 2D Contours. *Robotics and Autonomous Systems*, 2021, 142, pp.#103798. 10.1016/j.robot.2021.103798 . hal-02558064v3

**HAL Id: hal-02558064**

**<https://hal.science/hal-02558064v3>**

Submitted on 22 Apr 2021

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

## Highlights

### **2 Vision-based Manipulation of Deformable and Rigid Objects Using Subspace Projections of 2D Contours**

4 Jihong Zhu, David Navarro-Alarcon, Robin Passama, Andrea Cherubini

- 6 • We present a unique framework for manipulating both rigid and de-  
formable objects.
- Our framework is model-free and requires a short initialization phase.
- 8 • Our framework does not require camera calibration, and works with  
different camera poses.

# Vision-based Manipulation of Deformable and Rigid Objects Using Subspace Projections of 2D Contours

Jihong Zhu<sup>a,b</sup>, David Navarro-Alarcon<sup>c</sup>, Robin Passama<sup>a</sup>, Andrea  
Cherubini<sup>a</sup>

<sup>a</sup>*LIRMM - Université de Montpellier CNRS, 161 Rue Ada, 34090 Montpellier, France.*

`firstname.lastname@lirmm.fr`

<sup>b</sup>*Delft University of Technology and Honda Research Institute, Europe.*

`j.zhu-3@tudelft.nl`

<sup>c</sup>*The Hong Kong Polytechnic University, Department of Mechanical Engineering,  
Kowloon, Hong Kong. david.navarro-alarcon@polyu.edu.hk*

---

## Abstract

This paper proposes a unified vision-based manipulation framework using image contours of deformable/rigid objects. Instead of explicitly defining the features by geometries or functions, the robot automatically learns the visual features from processed vision data. Our method simultaneously generates—from the same data—both visual features and the interaction matrix that relates them to the robot control inputs. Extraction of the feature vector and control commands is done online and adaptively, and requires little data for initialization. Our method allows the robot to manipulate an object without knowing whether it is rigid or deformable. To validate our approach, we conduct numerical simulations and experiments with both deformable and rigid objects.

*Keywords:* Visual servoing, sensor-based control, deformable object manipulation.

---

## 1. Introduction

Humans are capable of manipulating both rigid and deformable objects. However, robotic researchers tend to consider the manipulation of these two classes of objects as separate problems. Unless otherwise mentioned, object rigidity is an implicit assumption in most manipulation tasks. On the other hand, methods designed for deformable object manipulation ([Sanchez et al.](#),

2018), are never applied on rigid objects. This paper presents our efforts in  
2 formulating a generalized framework for vision-based manipulation of both  
rigid and deformable objects, which does not require prior knowledge of the  
4 object’s mechanical properties.

In the visual servoing literature (Chaumette and Hutchinson, 2006), vec-  
6 tor  $\mathbf{s}$  denotes the set of features selected to represent the object in the image.  
These features represent both the object’s pose and its shape. We denote  
8 the process of selecting  $\mathbf{s}$  as *parameterization*. The aim of visual servoing is  
to minimize, through robot motion, the feedback error  $\mathbf{e} = \mathbf{s}^* - \mathbf{s}$  between  
10 the target  $\mathbf{s}^*$  and the current (i.e., measured) feature  $\mathbf{s}$ .

One of the initial works on vision-based manipulation of deformable ob-  
12 jects is presented in (Inoue, 1984) to solve a knotting problem by a topological  
model. Smith et al. developed a relative elasticity model, such that vision can  
14 be utilized without a physical model for the manipulation task (Smith et al.,  
1996). A classical model-free approach in manipulating deformable objects is  
16 developed in (Berenson, 2013). More recent research (Lagneau et al., 2020a)  
and (Lagneau et al., 2020b) proposes a method for online estimation of the  
18 deformation Jacobian, based on weighted least square minimization with a  
sliding window. In (Navarro-Alarcon et al., 2014) and (Navarro-Alarcon and  
20 Liu, 2018), the vision-based deformable objects manipulation is termed as  
*shape servoing*. An expository paper on the topic is available in (Navarro-  
22 Alarcon et al., 2019). A recent work on vision-based shape servoing of plastic  
material was presented in (Cherubini et al., 2020).

For a detailed survey on shape servoing we refer readers to (Sanchez  
24 et al., 2018). For *shape servoing*, commonly selected features are curvatures  
26 (Navarro-Alarcon et al., 2014), points (Wang et al.) and angles (Navarro-  
Alarcon and Liu, 2013). Laranjeira et al. proposed a catenary-based fea-  
28 ture for tethered management on wheeled and underwater robots (Laran-  
jeira et al., 2017, 2020). A more general feature vector is that containing the  
30 Fourier coefficients of the object contour (Navarro-Alarcon and Liu, 2018;  
Zhu et al., 2018). Yet, all these approaches require the user to specify a  
32 model, e.g., the object geometry (Wang et al.; Navarro-Alarcon and Liu,  
2013; Navarro-Alarcon et al., 2014) or a function (Laranjeira et al., 2017;  
34 Navarro-Alarcon and Liu, 2018; Zhu et al., 2018) for selecting the feature. Al-  
ternative data-driven (hence, model-free) approaches rely on machine learn-  
36 ing. Nair et al. combine learning and visual feedback to manipulate ropes in  
(Nair et al., 2017). Li et al. approximate the deformation and camera model  
38 using a neural network (Li et al., 2018). The authors of (Hu et al., 2019) em-

ploy deep neural networks to manipulate deformable objects given their 3D point cloud. All these methods rely on (deep) connectionist models, which invariably require training through an extensive data set. The collected data has to be diverse enough to generalize the model learnt by this type of networks. Instead of relying on algorithmic solutions, (She et al., 2020) utilizes a vision-based tactile sensor (GelSight) for manipulating cables.

It is noteworthy that some of the above mentioned methods may apply to rigid objects. Yet, none of the previous works has investigated the possibility of this extension nor reported its experimental validation, as we do in this paper.

The trend in visual servoing, when *controlling the pose of rigid objects* is to find features which are independent from the object characteristics. Following this trend, (Chaumette, 2004) proposes the use of image moments. More recently, researchers have proposed direct visual servoing (DVS) methods, which eliminate the need for user-defined features and for the related image processing procedures. The pioneer DVS works (Collewet et al., 2008; Collewet and Marchand, 2011) propose using the whole image luminance to control the robot, leading to “photometric” visual servoing. Bakthavatchalam et al. join the two ideas by introducing photometric moments (Bakthavatchalam et al., 2013). A subspace method (Marchand, 2019) can further enhance the convergence of photometric visual servoing, via Principal Component Analysis (PCA). This method was first introduced for visual servoing in (Nayar et al., 1996). In that work, using an eye-in-hand setup, the image was compressed to obtain a low-dimensional vector for controlling the robot to a target pose. Similarly, the authors of (Deguchi and Noguchi, 1996) transformed the image into a lower dimensional hyper surface, to control the robot position via in-hand camera feedback. However, DVS generally considers rigid and static scenes, where the robot controls the motion of the camera (eye-in-hand setup) to change only the image viewpoint, and not the environment. These constraints on the setup avoid breaking the Lambertian hypothesis that is needed, since DVS relies on the raw image luminance, which should not vary with the viewpoint. For this reason, to our knowledge, DVS was never applied to object manipulation, since changes in the pose and/or shape of the object would break the Lambertian assumption. This is not the case of feature-based methods (such as the one we present here), as long as the feature is chosen to be reliable even when the viewpoint and/or scene change.

Compared with the above-mentioned works, our paper presents the fol-

lowing original contributions:

1. We propose to use a feature vector – based on PCA of sampled 2D contours – for model-free manipulation of both deformable and rigid objects.
2. We exploit the linear properties of PCA and of the local interaction matrix, to initialize our algorithm with little data – the same data for feature vector extraction and for interaction matrix estimation.
3. We report experiments using the same framework to manipulate objects with different unknown geometric and mechanical properties.

The paper is organized as follows. Sect. 2 presents the problem. Sect. 3 outlines the framework. Sect. 4 elaborates on the methods. In Sect. 5, we analyze and verify the methods by numerical simulations. Then, Sect. 6 presents the robotic experiments and we conclude in Sect. 7.

## 2. Problem statement

In this work, we aim at solving object manipulation tasks with visual feedback. We rely on the following hypotheses:

- The shape and pose of the object are represented by its 2-D contour on the image as seen from a camera fixed in the robot workspace (eye-to-hand setup). We denote this contour as

$$\mathbf{c} = [\mathbf{p}_1 \ \cdots \ \mathbf{p}_K]^T \in \mathbb{R}^{2K}, \quad (1)$$

where  $\mathbf{p}_j = [u_j \ v_j] \in \mathbb{I}$  denotes the  $j$ th pixel of the contour in the image  $\mathbb{I}$ .

- The contour is always entirely visible in the scene and there are no occlusions.
- One of the robot’s end-effectors holds one point of the object (we consider the grasping problem to be already solved). At each control iteration  $i$ , its pose is  $\mathbf{r}_i \in \mathbb{SE}(3)$ , and it can execute motion commands  $\delta \mathbf{r}_i \in \mathbb{SE}(3)$  that drive the robot so that  $\mathbf{r}_{i+1} = \mathbf{r}_i + \delta \mathbf{r}_i$ .
- The target constant shape (i.e., contour) of the object,  $\mathbf{c}^*$ , is physically reachable with shaping motions of the grasping point  $\mathbf{r}$ . To ensure this hypothesis, one can first command the robot to verify that it can move the shape to  $\mathbf{c}^*$ .

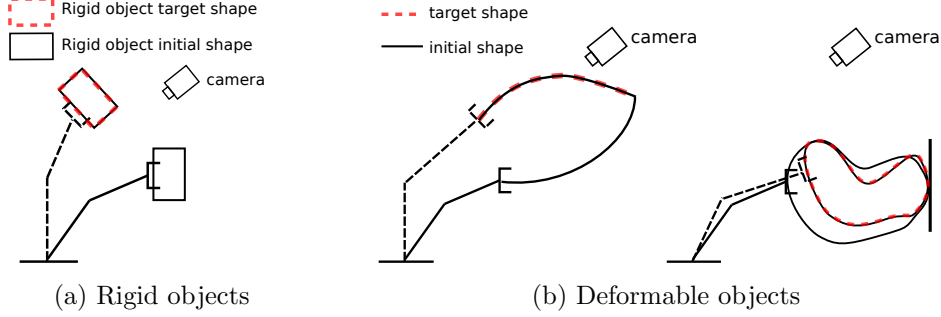


Figure 1: Vision-based manipulation of rigid and deformable objects. For rigid objects (left): control pose (translation and rotation). For deformable objects (right): control the pose, and also shape.

**Problem Statement.** Given a target shape of the object, represented by a constant contour vector  $\mathbf{c}^*$ , we aim at designing a vision-based controller that generates a sequence of robot motions  $\delta \mathbf{r}_i$  to drive the initial contour to the target one.

The controller should work without any knowledge of the object physical characteristics, i.e., for both rigid and deformable objects. In the latter case, we assume that the deformation is homogeneous. Since rigid and deformable objects behave differently during manipulation, we set the following manipulation goals:

- Rigid objects: move them to a target pose (see Fig. 1a).
- Deformable objects: move them to a target pose with a target shape (see Fig. 1b).

The formulation of the problem is general, but due to challenges in perception (discussed in Sect. 7), we carried out the cases of study with movements in  $\mathbb{SE}(2)$ .

### 3. Preliminary

In this section, we present an overview of the proposed approach, motivated by the problem analysis. Throughout the paper, we use  $\mathbf{c}$  to indicate the *object contour* and  $\mathbf{s}$  as the *feature vector* obtained from the contour.

The subscript  $i$  indicates the instance of the variable at iteration  $i$  (e.g.,  $\mathbf{c}_i$  is the contour at iteration  $i$ ).

We can work directly on the object shape space by selecting the contour as the feature vector  $\mathbf{s} \equiv \mathbf{c} \in \mathbb{R}^{2K}$ . With image and data processing, we can extract a fixed number of ordered (i.e., identified) contour points to represent the shape/pose of the object. However, this will result in an unnecessarily large dimension of the feature vector (e.g., if  $K = 50$ ,  $\mathbf{s}$  has 100 components). The high dimensional feature vector increases the computation demand and complicates the control due to the high under-actuation of the system. Therefore, instead of working on this feature vector, we work on one with smaller dimensions. To this end, we split the problem into two sub-problems: *parameterization* and *control*, see Fig. 2.

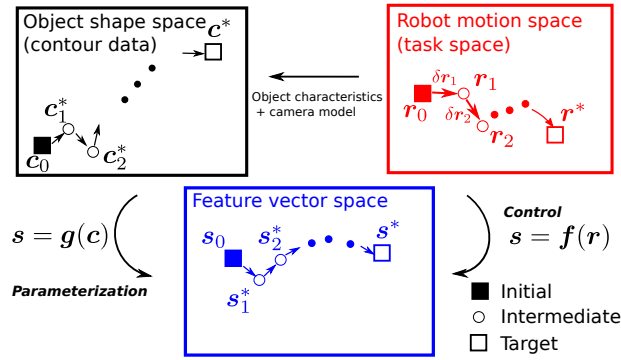


Figure 2: Graphic representation of the vision-based manipulation problem, with its two sub-problems, *parameterization* and *control*.

*Parameterization* consists in representing the contour via a compact feature vector  $\mathbf{s} \in \mathbb{R}^k$ , such that  $k \ll 2K$ . We denote this representation as  $\mathbf{s} = \mathbf{g}(\mathbf{c})$ . We introduce the method for parameterization in Sect. 4.1.

*Control* consists in computing robot motions  $\delta \mathbf{r}_1, \delta \mathbf{r}_2, \dots$ , so that the object's representation  $\mathbf{s}$  converges to the target  $\mathbf{s}^*$ . *Control* can be broken down to solving the optimization problem:

$$\mathbf{r}^* = \arg \min_{\mathbf{r}} (\mathbf{f}(\mathbf{r}) - \mathbf{s}^*) \quad (2)$$

where  $\mathbf{s} = \mathbf{f}(\mathbf{r})$  denotes the mapping between robot pose and feature vector, which is assumed to be smooth and generally nonlinear. The smoothness assumption requires that the objects' contour is at least twice differentiable



with respect to the robot motion. If we know the analytic solution to  $\mathbf{f}(\mathbf{r})$ , we can solve (2) and obtain the target shape in a single iteration by commanding  $\mathbf{r}^*$ .

A solution to this problem is to approximate the full mapping  $\mathbf{f}(\mathbf{r})$  from sensor observations. Classic deep learning-based approaches typically require a long training phase to collect vast and diverse data for approximating  $\mathbf{f}(\mathbf{r})$ . In some cases (for instance, robotics surgery), it is not possible to collect such data beforehand. Moreover, if the object changes, new data has to be collected to retrain the model, leading to a cumbersome process. In this paper instead, we aim at doing the data collection online, with minimum initialization.

Thus, instead of estimating the full nonlinear mapping  $\mathbf{f}(\mathbf{r})$ , we divide it into piece-wise linear models (Sang and Tao, 2012) at successive equilibrium points. The locality assumption refers to both the time and spatial dimensions. These models are considered time invariant in the neighbourhood of the equilibrium points. We then compute the control law for each linear model and apply it to the robot end-effector. We will dedicate Sections 4.2 and 4.3 to the local models and Sections 4.4 and 4.5 to derive the control inputs and to analyze (local) stability.

## 4. Methodology

Given a target shape  $\mathbf{c}^*$ , we define an intermediate local target  $\mathbf{c}_i^*$  at each  $i = 1, 2, \dots$  (see Fig. 2). At the  $i^{\text{th}}$  iteration, the robot autonomously generates a local mapping  $\mathbf{g}_i$  to produce the feature vector  $\mathbf{s}_i = \mathbf{g}_i(\mathbf{c}_i)$ . The robot then finds the local mapping  $\mathbf{s}_i = \mathbf{f}_i(\mathbf{r}_i)$  online.

Consider at the current time instant  $i$ , the shape  $\mathbf{c}_i$ , the intermediate target  $\mathbf{c}_i^*$  and the local parameterization  $\mathbf{g}_i$ . We can transform shape data into a feature vector by:

$$\mathbf{s}_i = \mathbf{g}_i(\mathbf{c}_i), \quad \mathbf{s}_i^* = \mathbf{g}_i(\mathbf{c}_i^*). \quad (3)$$

The linearized version of  $\mathbf{s} = \mathbf{f}(\mathbf{r})$  centered at  $(\mathbf{s}_i, \mathbf{r}_i)$  is then:

$$\delta \mathbf{s}_i = \mathbf{L}_i \delta \mathbf{r}_i, \quad (4)$$

with

$$\begin{aligned} \mathbf{L}_i &= \frac{\partial \mathbf{f}_i}{\partial \mathbf{r}} \big|_{\mathbf{r}=\mathbf{r}_i}, \\ \delta \mathbf{s}_i &= \mathbf{s}_{i+1} - \mathbf{s}_i, \\ \delta \mathbf{r}_i &= \mathbf{r}_{i+1} - \mathbf{r}_i. \end{aligned} \tag{5}$$

The matrix  $\mathbf{L}_i$  represents a local mapping, referred to as the interaction matrix in the visual servoing literature (Chaumette and Hutchinson, 2006). If  $\mathbf{L}_i$  can be estimated online at each iteration  $i$ , then, we can design one-step control laws to drive  $\mathbf{s}_i$  towards  $\mathbf{s}_i^*$ .

After the robot has executed the motion command  $\delta \mathbf{r}_i$ , we update the next target to be  $\mathbf{s}_{i+1}^*$ , and so on, until it reaches the final target  $\mathbf{s}^*$ . Although the validity region of this local mapping is smaller than that of the original nonlinear mapping, it enables to use an online training approach that requires less data and reduced computational demand.

Figure 3 shows the building blocks of the overall framework. In this section, we focus on the red dashed part of the diagram. We will elaborate on each red block in the subsequent subsections. The blue block represents the image processing pipeline that will be discussed in Sect. 6.1.

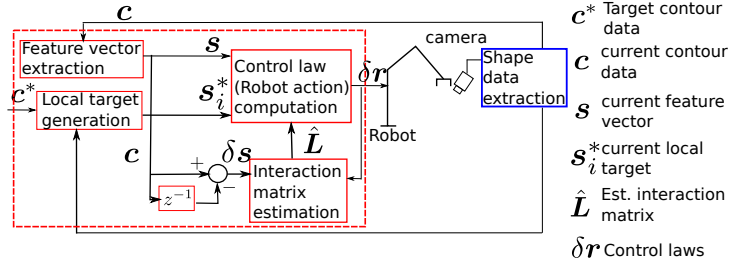


Figure 3: The block diagram that represents the overall framework.

#### 4.1. Feature vector extraction

There are many ways to parameterize  $\mathbf{c}$  in order to reduce its dimension. One of the prominent dimension reduction methods is Principal Component Analysis (PCA). PCA finds a new orthogonal basis for high-dimensional data. This enables projection of the data to lower dimension with the minimal sum of squared residuals. It was used in image processing (Zhang et al., 2010) and classification (Zeng et al., 2016). In visual servoing, the method was first introduced in (Nayar et al., 1996). PCA is proven to be an effective, yet easy to implement, algorithm for dimension reduction. By projecting to the new

orthogonal space, each feature component is linearly independent. Besides,  
 2 by checking the explained variance of a feature, we can intuitively measure  
 if it represents the original shape.

4 We apply PCA to reduce  $\mathbf{c} \in \mathbb{R}^{2K}$  to  $\mathbf{s} \in \mathbb{R}^k$ . To find the projection, we  
 collect  $M$  images with different shapes of the object and construct the data  
 6 matrix  $\mathbf{\Gamma} = [\mathbf{c}_1 \ \mathbf{c}_2 \ \cdots \ \mathbf{c}_M] \in \mathbb{R}^{2K \times M}$ . Then, we shift the columns of  $\mathbf{\Gamma}$  by  
 the mean contour  $\bar{\mathbf{c}} = \sum_{i=1}^M \mathbf{c}_i / M$ :

$$\bar{\mathbf{\Gamma}} = [\mathbf{c}_1 - \bar{\mathbf{c}} \ \mathbf{c}_2 - \bar{\mathbf{c}} \ \cdots \ \mathbf{c}_M - \bar{\mathbf{c}}] \in \mathbb{R}^{2K \times M}. \quad (6)$$

8 We then compute the covariance matrix  $\mathbf{C} = \bar{\mathbf{\Gamma}} \bar{\mathbf{\Gamma}}^T$ , and apply Singular Value  
 Decomposition (SVD) to it:

$$\mathbf{C} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T. \quad (7)$$

10 Once we have obtained the eigenvector matrix  $\mathbf{U} \in \mathbb{R}^{2K \times 2K}$ , we can move  
 on to select the first  $k$  columns<sup>1</sup> of  $\mathbf{U}$  denoted by  $\mathbf{U}(k) \in \mathbb{R}^{2K \times k}$ . Then,  
 12 the  $2K$ -dimensional contour  $\mathbf{c}$  can be projected into a smaller  $k$ -dimensional  
 feature vector  $\mathbf{s}$  as:

$$\mathbf{s} = \mathbf{U}^T(k)(\mathbf{c} - \bar{\mathbf{c}}) \in \mathbb{R}^k. \quad (8)$$

14 To assess the quality of this projection, we can compute the *explained*  
*variance* using the eigenvalue matrix  $\mathbf{\Sigma} \in \mathbb{R}^{2K \times 2K}$  in (7). By denoting the  
 16 diagonal entries of  $\mathbf{\Sigma}$  as  $\sigma_1, \dots, \sigma_{2K}$ , the explained variance of the first  $k$   
 components is:

$$\Upsilon(k) = \frac{\sum_{j=1}^k \sigma_j}{\sum_{j=1}^{2K} \sigma_j}. \quad (9)$$

18 where  $\Upsilon$  is a scalar between 0 and 1 (since  $\sigma_j > 0, \forall j$ ), indicating to what  
 extent the  $k$  components represent the original data (a larger  $\Upsilon$  suggests a  
 20 better representation).

Since PCA calculate features that lie on an orthonormal basis, these  
 22 features are linearly independent. For controlling  $n$  DoF, at least the same  
 number of independent visual features should be used. Therefore, we set  
 24  $k = n$  features.

---

<sup>1</sup>In the SVD algorithm, the first  $k$  columns correspond to the  $k$  largest eigenvalues of  
 matrix  $\Sigma$ .

#### 4.2. Local target generation

Let us now explain how we generate a local target contour  $\mathbf{c}_i^*$  given a current contour  $\mathbf{c}_i$  and final target contour  $\mathbf{c}^*$ . We also show this in Algorithm 1. The overall shape error is given by:

$$\mathbf{c}_e = \mathbf{c}^* - \mathbf{c}_i. \quad (10)$$

We define the intermediate target contour as:

$$\mathbf{c}_i^* = \mathbf{c}_i + \frac{1}{\eta} \mathbf{c}_e, \quad (11)$$

with  $\eta = 1, 2, \dots$  an integer that ensures that  $\mathbf{c}_i^*$  is a “good” local target for  $\mathbf{c}_i$  (i.e., the two are similar). Therefore, if we project the intermediate local data using the eigenvector matrix at the current iteration,  $\mathbf{U}_i \in \mathbb{R}^{2K \times 2K}$  (note that we are using the full projection matrix and not just the first  $k$  columns), the projection  $\mathbf{s}_i^p = \mathbf{U}_i(\mathbf{c}_i^* - \bar{\mathbf{c}}) \in \mathbb{R}^{2K}$  should fulfil:

$$\Psi(k) = \frac{\sum_{j=1}^k |s_{i,j}^p|}{\sum_{j=1}^{2K} |s_{i,j}^p|} \geq \epsilon, \quad (12)$$

with  $\epsilon \in [0; 1]$  a threshold and  $s_{i,j}^p$  the  $j$ -th component of the projection. Then, we select the first  $k$  components in  $\mathbf{s}_i^p$  to be the local target  $\mathbf{s}_i^* \in \mathbb{R}^k$ .

Algorithm 1 outlines the steps for computing the local intermediate targets, so that:

- they are near the final target,
- the corresponding feature vector can be extracted with the current learned projection matrix.

**Remark 1.** The reachability of a local target can only be verified with a global deformation model which we want to avoid identifying in our methods. We will further discuss this issue in the Conclusion (Sect. 7).

#### 4.3. Interaction matrix estimation

Let us consider the current contour  $\mathbf{c}_i$  and the local target  $\mathbf{c}_i^*$ . In this section, we show how we can implement the PCA and model estimation

---

**Algorithm 1** Local target generation
 

---

```

localTargetFound = false
 $\Psi_0 = 0$ 
 $\eta = 1$ 
while not localTargetFound do
   $\mathbf{c}_i^* = \mathbf{c}_i + \frac{1}{\eta}(\mathbf{c}^* - \mathbf{c}_i)$ 
   $\mathbf{s}_i^p = \mathbf{U}_i \mathbf{c}_i^*$ 
   $\Psi_\eta = \sum_{j=1}^k |s_{i,j}^p| / \sum_{j=1}^{2K} |s_{i,j}^p|$ 
  if  $\Psi_\eta \geq \epsilon$  or  $\Psi_\eta < \Psi_{\eta-1}$  then
    localTargetFound = true
     $\mathbf{s}_i^* = [\mathbf{I} \ 0] \mathbf{s}_i^p$ 
  end if
   $\eta = \eta + 1$ 
end while

```

---

together and online. We denote the robot motions and corresponding object contours over the last  $M$  iterations (prior to iteration  $i$ , with  $i \geq M$ ) as:

$$\begin{aligned} \Delta \mathbf{R}_i &= [\delta \mathbf{r}_{i-M+1} \ \delta \mathbf{r}_{i-M+2} \cdots \delta \mathbf{r}_i] \in \mathbb{R}^{n \times M} \\ \mathbf{\Gamma}_i &= [\mathbf{c}_{i-M} \ \mathbf{c}_{i-M+1} \ \mathbf{c}_{i-M+2} \cdots \mathbf{c}_i] \in \mathbb{R}^{2K \times (M+1)}, \end{aligned} \quad (13)$$

with  $M$  the number of data samples collected during initialization, i.e., the size of the sliding window used for model adaptation (see Sect. 4.5).

By selecting  $k = n$  (note that  $n$  is also the number of DoFs of the robot manipulator we considered in the task execution), we compute the projection matrix  $\mathbf{U}_i(n) \in \mathbb{R}^{2K \times n}$ , from  $\mathbf{\Gamma}_i$  and  $\bar{\mathbf{c}}_i$  via (6) and (7). Then, using  $\mathbf{U}_i(n)$ , we project current contour  $\mathbf{c}_i$ , target contour  $\mathbf{c}_i^*$  and shape matrix  $\mathbf{\Gamma}_i$ :

$$\begin{aligned} \mathbf{s}_i &= \mathbf{U}_i(n)^T (\mathbf{c}_i - \bar{\mathbf{c}}_i) \in \mathbb{R}^n, \\ \mathbf{s}_i^* &= \mathbf{U}_i(n)^T (\mathbf{c}_i^* - \bar{\mathbf{c}}_i) \in \mathbb{R}^n, \\ \mathbf{S}_i &= \mathbf{U}_i(n)^T \bar{\mathbf{\Gamma}}_i = [\mathbf{s}_{i-M} \ \mathbf{s}_{i-M+1} \ \cdots \ \mathbf{s}_i] \in \mathbb{R}^{n \times (M+1)}. \end{aligned} \quad (14)$$

In (14),  $\bar{\mathbf{\Gamma}}_i$  is normalized by  $\bar{\mathbf{c}}_i$  as in (6). We can then compute  $\Delta \mathbf{S}_i$  from (5) and (14), by subtracting consecutive columns of  $\mathbf{S}_i$ :

$$\Delta \mathbf{S}_i = [\delta \mathbf{s}_{i-M+1} \ \delta \mathbf{s}_{i-M+2} \cdots \delta \mathbf{s}_i] \in \mathbb{R}^{n \times M}. \quad (15)$$

Using  $\Delta \mathbf{S}_i \in \mathbb{R}^{n \times M}$  and  $\Delta \mathbf{R}_i \in \mathbb{R}^{n \times M}$  we can now estimate the local interaction matrix  $\mathbf{L}_i \in \mathbb{R}^{n \times n}$  at iteration  $i$ . We assume that near this iteration,

the system remains linear and time invariant:  $\mathbf{L}_i$  is constant. Using the local  
 2 linear model (4), we can write the following:

$$\Delta \mathbf{S}_i = \mathbf{L}_i \Delta \mathbf{R}_i. \quad (16)$$

Our goal then is to solve for  $\mathbf{L}_i$ , given  $\Delta \mathbf{S}_i$  and  $\Delta \mathbf{R}_i$ . Note that this is an  
 4 overdetermined linear system (with  $n \times M$  equations for  $n^2$  unknowns). Let us  
 consider  $\Delta \mathbf{R}_i \in \mathbb{R}^{n \times M}$  has full row rank. Note this sufficiently implies  $M \geq$   
 6  $n$ . With this prerequisite,  $\text{rank}(\Delta \mathbf{R}_i) = n$ . Therefore,  $\text{rank}(\Delta \mathbf{R}_i \Delta \mathbf{R}_i^T) = n$ ,  
 and its inverse exists. We post multiply (16) by  $\Delta \mathbf{R}_i^T$ :

$$\Delta \mathbf{S}_i \mathbf{R}_i^T = \mathbf{L}_i \Delta \mathbf{R}_i \Delta \mathbf{R}_i^T. \quad (17)$$

8 Then, since  $\Delta \mathbf{R}_i \Delta \mathbf{R}_i^T$  is invertible, the  $\mathbf{L}_i$  that best fulfills (16) is:

$$\hat{\mathbf{L}}_i = \Delta \mathbf{S}_i \Delta \mathbf{R}_i^T (\Delta \mathbf{R}_i \Delta \mathbf{R}_i^T)^{-1}. \quad (18)$$

If, in practice, the full row rank condition of  $\Delta \mathbf{R}_i$  is not satisfied,  $\text{rank}(\Delta \mathbf{R}_i \Delta \mathbf{R}_i^T) <$   
 10  $n$  and  $\Delta \mathbf{R}_i \Delta \mathbf{R}_i^T$  becomes singular. Then, instead of (18), we can use  
 Tikhonov regularization:

$$\hat{\mathbf{L}}_i = \Delta \mathbf{S}_i \Delta \mathbf{R}_i^T (\Delta \mathbf{R}_i \Delta \mathbf{R}_i^T + \lambda \mathbf{I})^{-1}, \quad (19)$$

12 with  $\lambda$  an arbitrary (generally small) scalar.

Practically, this implies that one or more inputs motions do not appear  
 14 in  $\Delta \mathbf{R}_i$ . Therefore, we cannot infer the relationship between these motions  
 and the resulting feature vector changes. In this case it is better to increase  
 16  $M$  and obtain more data, so that  $\Delta \mathbf{R}_i$  has full row rank.

Instead of computing the interaction matrix, it is also possible to directly  
 18 compute its inverse, since this guarantees better control properties (Lapresté  
 et al., 2004). With the same data, one can re-write (16) as:

$$\mathbf{L}_i^\oplus \Delta \mathbf{S}_i = \Delta \mathbf{R}_i. \quad (20)$$

20 We can also solve (20) with Tikhonov regularization:

$$\hat{\mathbf{L}}_i^\oplus = \Delta \mathbf{R}_i \Delta \mathbf{S}_i^T (\Delta \mathbf{S}_i \Delta \mathbf{S}_i^T + \lambda \mathbf{I})^{-1}. \quad (21)$$

#### 4.4. Control law and stability analysis

2 We can now control the robot, with either of the following strategies:

$$\delta \mathbf{r}_i = -\alpha \hat{\mathbf{L}}_i^\dagger (\mathbf{s}_i - \mathbf{s}_i^*), \quad (22)$$

if one estimates the interaction matrix with (19), where  $^\dagger$  denotes the pseudo-inverse, or:

$$\delta \mathbf{r}_i = -\alpha \hat{\mathbf{L}}_i^\oplus (\mathbf{s}_i - \mathbf{s}_i^*) \quad (23)$$

if one estimates the inverse of the interaction matrix with (21). In both equations,  $\alpha > 0$  is an arbitrary control gain.

**Proposition 1.** *Consider that locally, the model (4) closely approximates the interaction matrix  $\mathbf{L}_i = \hat{\mathbf{L}}_i$ . For  $M$  number of linearly independent displacement vectors  $\delta \mathbf{r}$  such that the interaction matrix  $\hat{\mathbf{L}}_i$  is invertible, the update rule (22) asymptotically minimizes the error  $\mathbf{e}_i = \mathbf{s}_i^* - \mathbf{s}_i$ , where  $\mathbf{s}_i^*$  is the local target.*

12 *Proof.* With  $\delta \mathbf{s}_i = \mathbf{s}_{i+1} - \mathbf{s}_i$ , we can write (4) in discretized form as

$$\mathbf{s}_{i+1} = \mathbf{s}_i + \mathbf{L}_i \delta \mathbf{r}_i. \quad (24)$$

From the definition of  $\mathbf{e}_i$  we have (Note here the target  $\mathbf{s}_i^*$  is not updated with  $i$  since we want to prove local convergence to a constant target):

$$\begin{aligned} \mathbf{e}_i &= \mathbf{s}_i^* - \mathbf{s}_i \\ \mathbf{e}_{i+1} &= \mathbf{s}_i^* - \mathbf{s}_{i+1} \end{aligned} \quad (25)$$

Taking (24) into (25):

$$\begin{aligned} \mathbf{e}_{i+1} &= \mathbf{s}_i^* - \mathbf{s}_{i+1} \\ &= \mathbf{s}_i^* - \mathbf{s}_i - \mathbf{L}_i \delta \mathbf{r}_i \\ &= \mathbf{e}_i - \mathbf{L}_i \delta \mathbf{r}_i. \end{aligned} \quad (26)$$

We replace  $\delta \mathbf{r}_i$  in (26) with the control (22), the error dynamic is then:

$$\begin{aligned} \mathbf{e}_{i+1} &= \mathbf{e}_i - \alpha \mathbf{L}_i \mathbf{L}_i^{-1} (\mathbf{s}_i^* - \mathbf{s}_i) \\ &= \mathbf{e}_i - \alpha \mathbf{e}_i = (1 - \alpha) \mathbf{e}_i. \end{aligned} \quad (27)$$

is asymptotically stable for  $\alpha \in [0; 1]$ . This can be proved by considering the Lyapunov function

$$\mathcal{V}(\mathbf{e}) = \mathbf{e}^T \mathbf{e}. \quad (28)$$

Using the error dynamic (27), one can derive:

$$\begin{aligned}\Delta\mathcal{V} &= \mathcal{V}(\mathbf{e}_{i+1}) - \mathcal{V}(\mathbf{e}_i) \\ &= \mathbf{e}_i^T((1-\alpha)^2 - 1)\mathbf{e}_i < 0.\end{aligned}\tag{29}$$

This proves the local asymptotic stability of the error  $\mathbf{e}$  using our inputs. ■

#### 2 4.5. Model adaptation

Since both the projection matrix  $\mathbf{U}^T(n)$  and the interaction matrix are  
4 local approximations of the full nonlinear mapping, they need to be updated  
constantly. We choose a receding window approach with window size  $M$ .

6 At current iteration  $i$ , we estimate the projection matrix  $\mathbf{U}_i^T$  and local  
interaction matrix  $\mathbf{L}_i$  with  $M$  samples of the most recent data. Using the  
8 interaction matrix and the local target  $\mathbf{c}_i^*$ , we can derive the one-step com-  
mand  $\delta\mathbf{r}_i$  by (22). Once we execute the motion  $\delta\mathbf{r}_i$ , a new contour data  $\mathbf{c}_{i+1}$   
10 is obtained. We move to the next iteration  $i + 1$ . A new pair of input and  
shape data  $[\delta\mathbf{r}_i, \mathbf{c}_{i+1}]$  is obtained. We shift the window by deleting the oldest  
12 data in the window and add in the new data pair. Then, using the shifted  
window, we compute one step control at iteration  $i + 1$ .

14 The receding window approach ensures that, at each iteration, we are us-  
ing the latest data to estimate the interaction matrix. The overall algorithm  
16 is initialized with small random motions around the initial configuration.  
First,  $M$  samples of shape data and the corresponding robot motions are  
18 collected. With this initialization, we can simultaneously solve for the pro-  
jection matrix and estimate the initial interaction matrix using the methods  
20 described in Sect. 4.1 and 4.3. Using the projection matrix and the ini-  
tial/target shapes, we can then find an intermediate target (see Sect. 4.2).

22 We consider quasi-static deformation. Hence, at each iteration the sys-  
tem is in equilibrium and can be linearized according to (4). The data that  
24 best captures the current system are the most recent ones. The choice of  $M$   
is a trade-off between locality and richness. For fast varying deformations<sup>2</sup>,  
26 we would expect to reduce  $M$  since a larger  $M$  will hinder the locality as-  
sumption. Yet, if  $M$  is too small, it affects the estimation of  $\hat{\mathbf{L}}_i$  (refer to the  
28 detailed discussion in Sect. 4.3).

---

<sup>2</sup>The notion of fast or slow varying depends on both the speed of manipulation, and  
on the objects deformation characteristics (which affect the rate of change in shapes) with  
regard to the image processing time.



## 5. Simulation results

In this section, we present the numerical simulations that we ran to validate our method.

### 5.1. Simulating the objects

We ran simulations on MATLAB (R2018b) with two types of objects: a rigid box and a deformable cable, both constrained to move on a plane. The rigid object is represented by a uniformly sampled rectangular contour. The controllable inputs are its position and orientation. For the cable, we developed a simulator, which is publicly available at <https://github.com/Jihong-Zhu/cableModelling2D>. The simulator relies on the differential geometry cable model introduced in (Wakamatsu and Hirai, 2004), with the shape defined by solving a constrained optimization problem. The underlying principle is that the object’s potential energy is minimal for the object’s static shape (Wakamatsu et al., 1995). Position and orientation constraints (imposed at the cable ends) are input to the simulator. The output is the sampled cable. Figures 4 – 6, 9, 10 show simulated shapes of cables and rigid boxes. We choose  $K = 50$  samples for both rigid objects and cables. The camera perspective projection is simulated, with optical axis perpendicular to the plane.

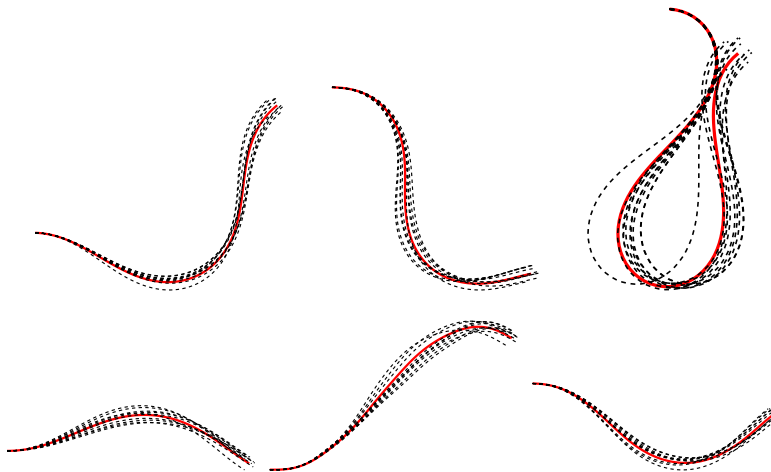


Figure 4: Six trials conducted to test various choices of feature dimension  $k$  for a cable. In each sub-figure, the solid red lines are the initial shapes and the dashed black are the shapes resulting from 10 random motions of the right tip (translations limited to  $\pm 5\%$  of the length, rotations limited to  $\pm 5^\circ$ ).

## 5.2. Selecting the feature dimension $k$

To check whether choosing  $k = n$  can represent the shape accurately, we simulate 6 trials with distinct initial shapes of a cable. The dimension of the robot motion vector  $\delta \mathbf{r}$  is  $n = 3$  (two translations and one rotation of the right tip), and the motions are limited: each translation to  $\pm 5\%$  of the cable length and the rotation to  $\pm 5^\circ$ . This range of motion gives a rule of thumb, which we have used for generating random movements throughout our experiments. For each trial, we command  $M = 10$  random motions around the initial shape using our simulator. Figure 4 shows the 6 initial cable shapes (solid red) and the resulting shapes from 10 random movements (dashed black).

For each trial, we apply PCA to map the cable contour  $\mathbf{c} \in \mathbb{R}^{2K}$  to feature vector  $\mathbf{s} \in \mathbb{R}^k$ , as explained in Sect. 4.1. We do this for  $k = 1, 2$  and 3 and for each of these 18 experiments, we calculate the explained variance  $\Upsilon(k)$  with (9). Table 1 shows these explained variances. In all 6 trials,  $k = n = 3$  yields explained variances very close to 1. This result confirms that choosing  $k = n$  as the dimension of the feature vector gives an excellent representation of the shape data. It is also possible to select  $k = 2$ , since the first two components can represent more than 99% of the variance. Nevertheless, the simulation is noise-free. Therefore, although  $\Upsilon(k)$  increases little from  $k = 2$  to  $k = 3$ , this increase is not related to noise but to an actual gain in data information.

Table 1: Explained variance  $\Upsilon(k)$  for the 6 trials with small motion.

	trial 1	trial 2	trial 3	trial 4	trial 5	trial 6
$k = 1$	0.727	0.795	0.871	0.847	0.847	0.705
$k = 2$	0.992	0.995	0.996	0.997	0.997	0.994
$k = 3$	0.999	0.999	0.999	0.999	0.999	0.999

At this stage, it is legitimate to ask: *how does this scale to larger movements?* Figure 5 illustrates 10 cable shapes generated by large movements (angle variation:  $[-\frac{\pi}{2}, \frac{\pi}{2}]$ , maximum translation: 106%). Again, we apply PCA ( $M = 10$ ); Table 2 shows the  $\Upsilon(k)$  resulting from various values of  $k$ .

Table 2: Explained variance  $\Upsilon(k)$  computed with large motion.

k	0	1	2	3	4	5
$\Upsilon(k)$	0	0.5444	0.7218	0.8927	0.9919	0.9990

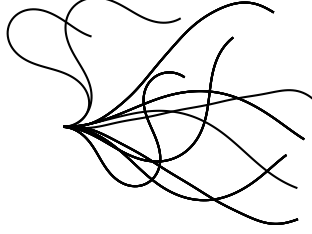


Figure 5: Ten distinctive cable shapes generated by large motion: angle variation:  $[-\frac{\pi}{2}, \frac{\pi}{2}]$ , maximum translation: 106% of the cable length.

Comparing Tables 1 and 2, it is noteworthy that  $\Upsilon(4)$  with large motion is smaller than  $\Upsilon(2)$  with small motion. There are two possible explanation here. One is that when shapes stays local, the local linear mapping  $\mathbf{L}$  in (4) remains constant and we need less features to characterize it; the more the shape varies, the more features we need. Another possible explanation is that for larger motions,  $M = 10$  shapes may be insufficient for PCA. Likely, the larger the changes, the larger the number of shapes  $M$  needed.

### 5.3. Manipulation of deformable objects

With our cable simulator, we can now test the controller to modify the shape from an initial to a target one. Again, the left tip of the cable is fixed, and we control the right tip with  $n = 3$  degrees of freedom (two translations and one rotation). Using the methods described in Sect. 4, we choose window size  $M = 5$ , the Tikhonov factor  $\lambda = 0.01$ , the local target threshold  $\epsilon = 0.8$ , the control gain  $\alpha = 0.01$ . To quantify the effectiveness of our algorithms in driving the contour to  $\mathbf{c}^*$ , we define a scalar measure: the Average Sample Error (ASE). At iteration  $i$ , with current contour  $\mathbf{c}_i$  it is:

$$\text{ASE} = \frac{\|\mathbf{c}_i - \mathbf{c}^*\|_2}{2K}. \quad (30)$$

A small ASE indicates that the current contour is near the target one. In Sect. 4.4, we have proved that our controller asymptotically stabilizes the feature vector,  $\mathbf{s}$  to  $\mathbf{s}^*$ . Hence, since we have also shown that  $\mathbf{s}$  is a “very good” representation” of  $\mathbf{c}$ , we also expect our controller to drive  $\mathbf{c}$  to  $\mathbf{c}^*$ , thus ASE to 0. This measure is also used in the real experiments.

Using the cable simulator, we compare the convergence of two control laws proposed in our paper (22) and (23) against a baseline algorithm in (Zhu et al., 2018) which uses Fourier parameters as feature. To make methods

compatible, we choose first order Fourier approximation. Note that this results in a feature vector of dimension of 6 (see (Zhu et al., 2018)) which is still twice the number  $k$  used in our method. We also normalize the computed control action and then multiply by the same gain factor 0.01.

We also introduced artificial noise to the contour data, to test the robustness of our method. For a unit length cable, we add Gaussian noise of zero mean and 0.01 standard deviation to the contour sample points. Fig. 6b shows that our algorithm converges in these conditions as well. It is worth mentioning that in robotic experiments, as shape data is obtained and sampled from the images, signal noise is inevitable. Yet, our framework is robust enough to still converge to the target shape/pose (see Sect. 6 for detailed real robot experiments).

Figure 6 shows two other simulation results: on the left, a reachable target and on the right an unreachable one. In Fig. 6a, the cable shape successfully evolves towards the target thanks to our controller (23). Figure 6c shows the starting shape (blue), unreachable target (dash black) and final shape (solid black) obtained using (23). Note that the controller gets stuck in a local minimum.

Figure 7 compares the evolution of ASE with our methods against the Fourier-based method for the reachable target; in the same figure, we also plot the evolution of ASE using (23) for the unreachable target. We can observe that our method provides faster convergence using half the features than (Zhu et al., 2018). Also, directly computing the inverse (23) provides faster convergence than (22). It is noteworthy to point out that the Fourier-based method requires a different parameterization for closed and open contours (see (Navarro-Alarcon and Liu, 2018) and (Zhu et al., 2018)), whereas in our framework, the parameterization can be kept the same. Last but not least, our approach is the only one among the three, which has been validated on both rigid and deformable objects.

#### 5.4. Comparison with the Broyden update law

The Broyden update law (Broyden, 1965), has been used to update the interaction matrix in classic visual servoing (Hosoda and Asada, 1994; Jagersand et al., 1997; Chaumette and Hutchinson, 2007) and shape servoing (Navarro-Alarcon et al., 2013).

In this section, we compare it with our method for updating the interaction matrix (19), which relies on a receding horizon. We will hereby show why the Broyden update law is not applicable in our framework.

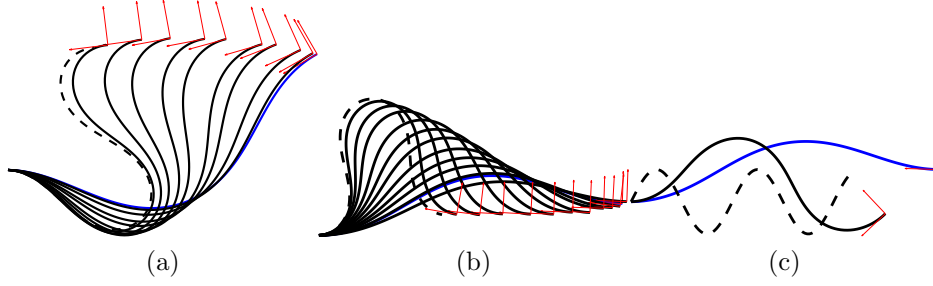


Figure 6: Cable manipulation with a single end-effector, moving the right tip. (a): a reachable target, the blue and black lines are the initial and intermediate shapes, respectively, and the dashed black line is the target shape. The red frame indicates the end-effector position and orientation generated by our controller. (b): Adding Gaussian noise with zero mean and 0.01 standard deviation to the shape data with a reachable target. (c): An unreachable target and the final shape obtained with our controller.

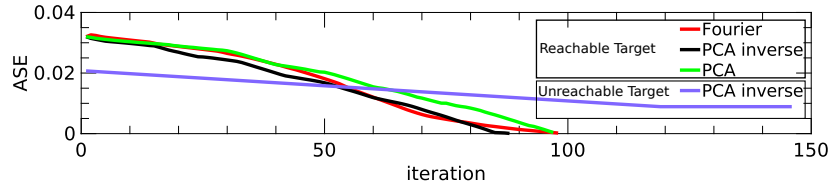


Figure 7: The evolution of the ASE of the simulated cable manipulation using our method against the Fourier-based method as baseline and the ASE of the unreachable target with (23).

The Broyden update is an iterative method for estimating  $\mathbf{L}_i$  at iteration  
 2 *i*. Its standard discrete-time formulation is:

$$\hat{\mathbf{L}}_i = \hat{\mathbf{L}}_{i-1} + \beta \frac{\delta \mathbf{s}_{i-1} - \hat{\mathbf{L}}_{i-1} \delta \mathbf{r}_{i-1}}{\delta \mathbf{r}_{i-1}^T \delta \mathbf{r}_{i-1}} \delta \mathbf{r}_{i-1}^T, \quad \forall \mathbf{r}_{i-1} \neq \mathbf{0} \quad (31)$$

with  $\beta \in [0; 1]$  an adjustable gain. Using our simulator, we estimate the  
 4 interaction matrix using both Broyden update (with three different values  
 of  $\beta$ ) and our receding horizon method (19). We then compare (with  $\hat{\mathbf{L}}$   
 6 estimated with either method) the one-step prediction of the resulting feature  
 vector:

$$\hat{\mathbf{s}}_{i+1} = \hat{\mathbf{L}}_i \mathbf{r}_i + \mathbf{s}_i, \quad (32)$$

8 with the ground truth  $\mathbf{s}_{i+1}$  from the simulator. The results (plotted in Fig. 8)  
 show that receding horizon outperforms all three Broyden trials. One possible  
 10 reason is that the components of  $\mathbf{s}$  fluctuate since (at each iteration) a new  
 matrix  $\mathbf{U}$  is used. These variations cause the Broyden method to accumulate  
 12 the result from old interaction matrices, and therefore perform badly on a  
 long term. This result contrasts with that of (Navarro-Alarcon et al., 2013),  
 14 where the Broyden method performs well since there is a fixed mapping from  
 contour data to feature vector. Another advantage of the receding horizon  
 16 approach is that it does not require any gain tuning.

### 5.5. Manipulation of rigid objects

18 The same framework can also be applied to rigid object manipulation.  
 Consider the problem of moving a rigid object to a certain position and  
 20 orientation via visual feedback. This time, the shape of the object does not  
 change, but its pose will (it can translate and rotate). We use the same  $M$ ,  
 22  $\lambda$ ,  $\epsilon$  and  $\alpha$  as for cable manipulation. We compare the convergence of two  
 control laws proposed in our paper (22) and (23) against a baseline using  
 24 image moments (Chaumette, 2004). The translation and orientation can be  
 represented with image moments and the analytic interaction matrix can be  
 26 computed as explained in (Chaumette, 2004)). To make methods compatible,  
 we normalize the computed control and then multiply it by the same factor  
 28 0.01.

Figure 9 shows two simulations where our controller successfully moves  
 30 a rigid object from an initial (blue) to a target (dashed black) pose using  
 control law (23). Figure 11 compares convergence of our methods against  
 32 the image moments method. We can observe that our method provides a

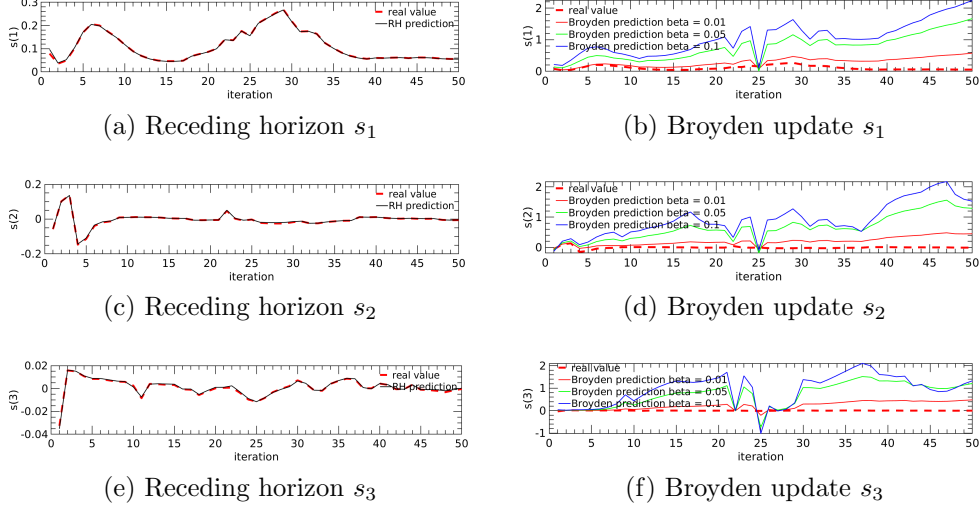


Figure 8: Comparison – for estimating  $\mathbf{s}$  – of the receding horizon approach (RH, left) and of the Broyden update (right, with three values of  $\beta$ ). The topmost, middle and bottom plots show the one step prediction of  $s_1$ ,  $s_2$  and  $s_3$ , respectively. In all plots, the dashed red curve is the ground truth from the simulator. The plots clearly show that the receding horizon approach outperforms all three Broyden trials.

slightly slower convergence. Directly computing the inverse (23) provides a  
 2 convergence similar to (22). Later, we will show why our method is slower.  
 Yet, the fact that it can be applied on both deformable and rigid objects  
 4 makes it stand out over the other techniques.

### 5.6. Feature analysis for rigid objects

6 In this section, we analyze locally what each component of the feature  
 vector represents, in the case of rigid object manipulation. To this end, we  
 8 apply  $M = 10$  random movements (rotation range  $[-0.11, 0.09]$ , maximum  
 translation 15% of the width) to multiple rigid rectangular objects (see Fig.  
 10 10). We compute the projection matrix as explained in Sect. 4.1, and trans-  
 form the contour samples to feature vectors. Then, we seek the relationship –  
 12 at each iteration – between the object pose  $x$ ,  $y$ ,  $\theta$  and the components of the  
 feature vector  $\mathbf{s}$  generated by PCA. To this end, we use bivariate correlation  
 14 (Feller, 2008) defined by:

$$\rho = \frac{E[(\xi - \bar{\xi})(\zeta - \bar{\zeta})]}{\sigma_{\xi}\sigma_{\zeta}}, \quad (33)$$

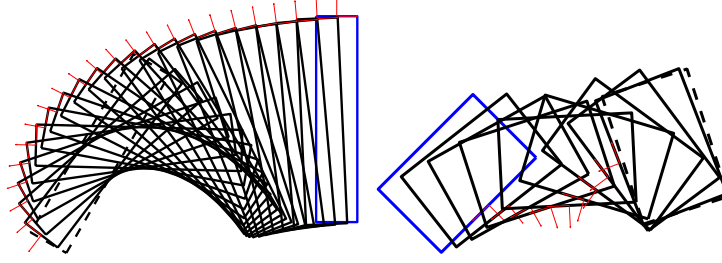


Figure 9: Manipulation of a rigid object with a single end-effector (red frame). The initial, intermediate and target contours are respectively blue, solid black and dashed black. Note that in both cases, our controller moves the object to the target pose.

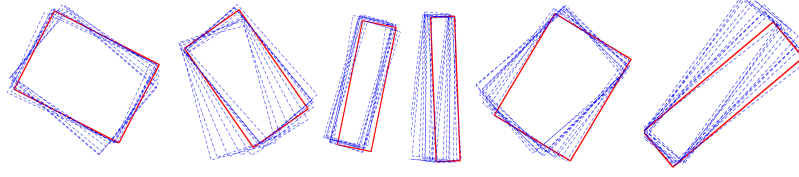


Figure 10: From an initial (red) pose, we generate 10 (dashed blue) random motions of a rigid object. This figure shows multiple examples of different rectangular rigid objects.

where  $\xi$  and  $\zeta$  are two variables with expected values  $\bar{\xi}$  and  $\bar{\zeta}$  and standard deviations  $\sigma_{\xi}$  and  $\sigma_{\zeta}$ . An absolute correlation  $|\rho|$  close to 1 indicates that the variables are highly correlated. All the simulations in Fig. 10 exhibit similar correlation between the computed feature vector and the object pose. In Table 3, we show one instance (Left first simulation in Fig. 10) of the correlation between variables, with high absolute correlations marked in red. It is clear from the table that each component in the feature vector relates strongly to one pose parameter. We further demonstrate the correlation in Fig. 12, where we plot the evolution of object poses and feature components. Note that  $s_2$  and  $\theta$  are negatively correlated. The slower convergence could be a result of the fact that, in contrast with image moments, here the extracted features and object pose are not completely decoupled. Yet, the main contribution of our method is that it can be directly used for both rigid and

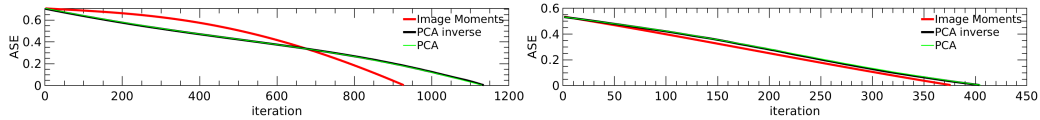


Figure 11: Evolution of ASE of the simulated rigid object manipulation using our method against image moments. Left: simulation in Fig. 9, Right: simulation in Fig. 9.



deformable objects. Therefore, it can be expected to be slower than methods  
 2 specifically designed for rigid objects (such as image moments).

Table 3: Correlation  $\rho$  between  $s_1, s_2, s_3$  and  $x, y, \theta$ .

	$x$	$y$	$\theta$
$s_1$	-0.2819	-0.3343	0.9887
$s_2$	0.2607	-0.8547	-0.0465
$s_3$	0.9230	0.3629	-0.1426

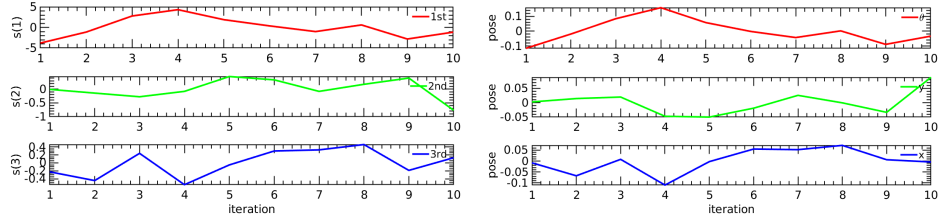


Figure 12: Progression of the auto-generated feature components (row 1, 3, 5:  $s_1, s_2, s_3$ ) vs. object pose (row 2, 4, 6:  $x, y, \theta$ ). We have purposely arranged the variables with high correlation with the same color.

## 6. Experiments

Figure 13 outlines our experimental setup. We use a KUKA LWR IV  
 4 arm. We constrain it to planar ( $n = 3$ ) motions  $\delta \mathbf{r}$ , defined in its base frame  
 6 (red in the figure): two translations  $\delta x$  and  $\delta y$  and one counterclockwise  
 8 rotation  $\delta \theta$  around  $z$ . A Microsoft Kinect V2 observes the object<sup>3</sup>. A Linux-  
 based 64-bit PC processes the image at 30 fps. In the following sections, we  
 10 first introduce the image processing for contour extraction, then present the  
 experiments.

### 6.1. Image processing

12 This section explains how we extract and sample the object contours  
 from an image. We have developed two pipelines, according to the kind of  
 14 contours (See Fig. 14): open (e.g., representing a cable) and closed. We  
 hereby describe the two.

<sup>3</sup>We only use the RGB image – not the depth – from the sensor.

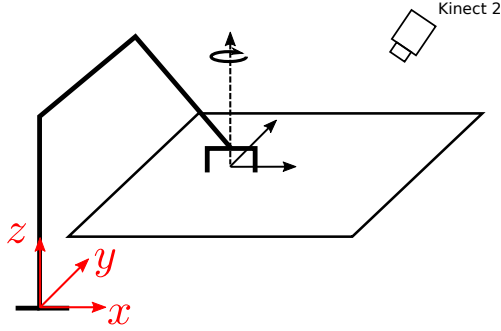


Figure 13: Overview of the experimental setup.

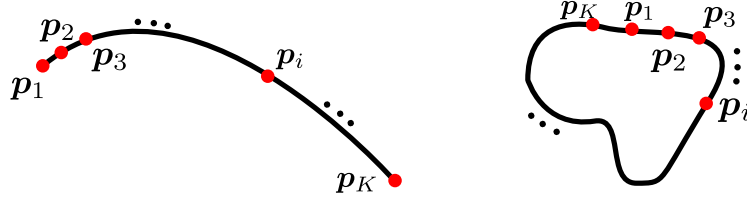


Figure 14: Open (left) and closed (right) contours can be both represented by a sequence of sample pixels in the image.

#### 6.1.1. Open contours

2 The overall pipeline for extracting an open contour is illustrated in Fig. 15  
 and **Algorithm 2**. On the initial image, the user manually selects a Region  
 4 of Interest (ROI, see Fig. 15a) containing the object. In this ROI, we apply  
 thresholding, followed by a morphological opening, to obtain a binary image  
 6 as in Fig. 15b. This image is dilated to generate a mask (Fig. 15c) used  
 to compute the new ROI for the following image. Figure 15e is the object  
 8 after a small manipulation motion and 15f shows the mask (in grey color)  
 which contains the cable. The OpenCV *findContour* function is applied to  
 10 binary image, then two contours are extracted based on the two known ends  
 of the cable, both are re-sampled (with same value of  $K$ ) using **Algorithm 2**  
 12 and finally merged (by interpolation, for each sample, between the two con-  
 tours' corresponding point) into the uniformly sampled open contour  $\mathbf{c}$  (see  
 14 Fig. 15d, where the green box indicates the end-effector).

#### 6.1.2. Closed contours

16 The procedure is shown in Fig. 16. For an object with uniform color  
 (in the experiment blue), we apply HSV segmentation, followed by Gaussian  
 18 blur of size 3, and finally the OpenCV *findContour* function, to get the object

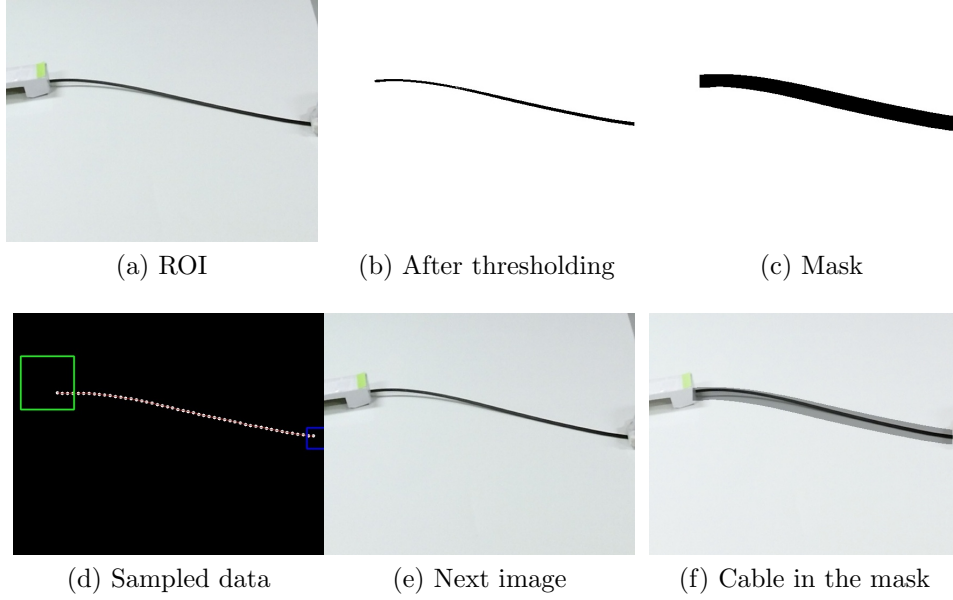


Figure 15: Image processing steps needed to obtain the sampled open contour of an object (here, a cable).

contour. The contour is then re-sampled using **Algorithm 2**. The starting point and the order of the samples is determined by tracking the grasping point (red dot in Fig. 16d) and the centroid of the object (blue dot). We obtain the vector connecting the grasping point to the centroid. Then, the starting sample is the one closest to this vector, and we proceed along the contour clockwise. Therefore, the extracted contour is ordered clockwise and always starting from the same point. This solves the contour ambiguity for symmetrical objects.

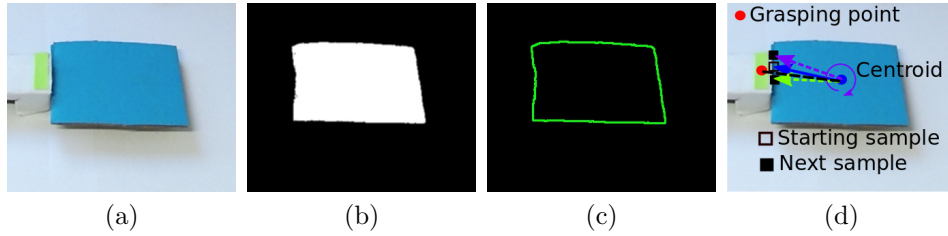


Figure 16: Image processing for getting a sampled closed contour: (a) original image, (b) image after thresholding and Gaussian blur, (c) extracted contour, (d) finding the starting sample and the order of the samples.

---

**Algorithm 2** Generate fixed number of points with uniform spacing

*Input*  $\mathbf{P}_I = [\mathbf{p}_I(1), \dots, \mathbf{p}_I(N)]$ : original ordered sampled data

*Input*  $K$  = target number of data samples

*Input*  $\epsilon$  = infinitesimal threshold for equality of distances

*Output*  $\mathbf{P}_O = [\mathbf{p}_O(1), \dots, \mathbf{p}_O(K)]$ : re-sampled data with uniform spacing.

---

```

1: compute the full length  $\mathcal{L}$  of  $\mathbf{P}_I$ .
2: compute desired distance per sample:  $\mu = \mathcal{L}/K$ 
3:  $l = 1, dist = 0$ 
4:  $\mathbf{p}_{curr} = \mathbf{p}_O(l) = \mathbf{p}_I(l)$ 
5:  $h = l = l + 1$ 
6: while  $l \leq N$  do
7:    $\mathbf{p}_{next} = \mathbf{p}_I(l)$ 
8:    $d = \|\mathbf{p}_{next} - \mathbf{p}_{curr}\|_2$ 
9:   if  $d + dist \leq \mu$  then
10:     $dist = dist + d$ 
11:     $\mathbf{p}_{curr} = \mathbf{p}_{next}$ 
12:     $l = l + 1$ 
13:   else
14:     $\mathbf{p}_{curr} = \mathbf{p}_O(h) = \mathbf{p}_{curr} + (\mathbf{p}_{next} - \mathbf{p}_{curr}) \frac{\mu - dist}{d}$ 
15:     $h = h + 1$ 
16:     $dist = 0$ 
17:   end if
18: end while
19: if  $|\mu - dist| < \epsilon$  then
20:    $\mathbf{p}_O(h) = \mathbf{p}_{curr}$ 
21: end if

```

---

## 6.2. Vision-based manipulation

In this section, we present the experiments that we ran to validate our algorithms, also visible at <https://youtu.be/gYfO2ZxZ5KQ>. To demonstrate the generality of our framework, we tested it with:

- Rigid objects represented by closed contours,
- Deformable objects represented by open contours (cables),
- Deformable objects represented by closed contours (sponges).

We carried out different experiments with a variety of initial and target contours and camera-to-object relative poses. The variety of both geometric and physical properties demonstrates the robustness of our framework. The variety of camera-to-object relative poses shows that—as usual in image-based visual servoing (Chaumette and Hutchinson, 2006)—camera calibration is unnecessary. The algorithm and parameters are the same in all experiments; the only differences are in the image processing, depending on the type of contour (closed or open, see Sect. 6.1).

We obtain the target contours by commanding the robot with predefined motions. Once the target contour is acquired, the robot goes back to the initial position, and then should autonomously reproduce the target contour. Again, we set the number of features  $k = n = 3$ , and use  $K = 50$  samples to represent the contour  $\mathbf{c}$ . We set the window size  $M = 5$ , both for obtaining the feature vector  $\mathbf{s}$  and the interaction matrix  $\mathbf{L}$ . We choose the control gain to be 0.01. A larger control gain may result in faster convergence, but could also lead to oscillation. The local target threshold  $\epsilon$  is set to 0.8. A higher threshold will result in closer local target and vice versa. The Tikhonov factor used to ensure numerical stability for matrix inversion, is set to  $\lambda = 0.01$ .

At the beginning of each experiment, the robot executes 5 steps of small<sup>4</sup> random motions to obtain the initial features and interaction matrix.

For all the experiments, we set the same termination condition at iteration  $i + 1$  using ASE defined in (30) such that:

1.  $\text{ASE}_i < 1$  pixel and
2.  $\text{ASE}_{i+1} \geq \text{ASE}_i$ .

---

<sup>4</sup>The notion of small is relative, and usually dependent on the size of the object the robot is manipulating. Refer to Sect. 5.2 (especially Fig. 4) for a discussion on this.

In the graphs that follow, we show the evolution of ASE in blue before the termination condition, and in red after the condition (until manual stop by the operator).

Figure 17 presents 8 experiments, one per column. Columns 1 – 3, 4 – 6 and 7 – 8 show respectively manipulation of: cable, rigid object and sponge. The first row presents the full RGB image obtained from Kinect V2. The second and third rows zoom in on the manipulation at the initial and final iterations. We track the end-effector in the image with a green marker for contour sampling. The target and current contours are drawn in red and blue, respectively.

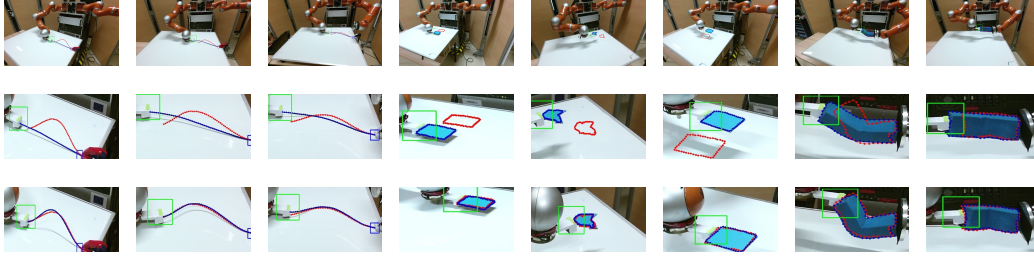


Figure 17: Eight experiments with the robot manipulating different objects. From left to right: a cable (columns 1 – 3), a rigid object (columns 4 – 6) and a sponge (columns 7 and 8). The first row shows the full Kinect V2 view, and the second and the third columns zoom in to show the manipulation process at the first and last iterations. The red contour is the target one, whereas the blue contour is the current one. The green square indicates the end-effector.

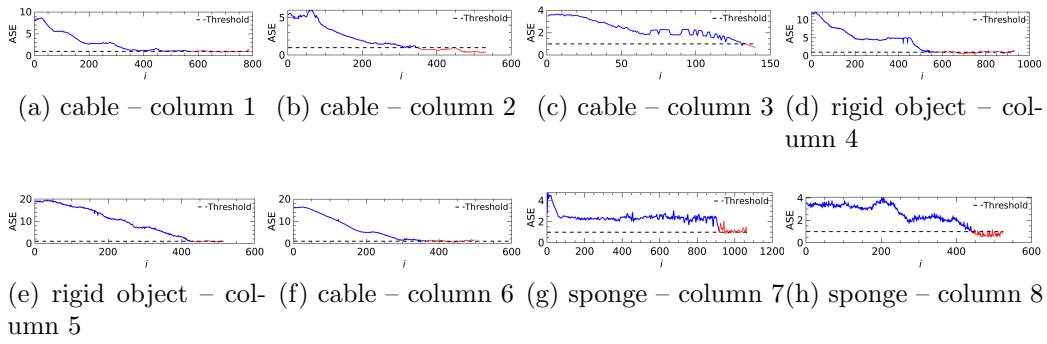


Figure 18: Evolution of  $e_i$  at each iteration  $i$ , for the 8 experiments of Fig. 17. The black dashed lines indicate the threshold  $ASE = 1$  pixel. The blue curves show  $e_i$  until the termination condition, whereas the red curves show the error until manual termination by the human operator.

Figure 18 shows the decreasing trend of error ASE for each experiment. The initial increase of ASE in the experiments can be due to the random motion at the beginning of the experiments. In general, we found that ASE is more noisy for the closed than for the open contour. This discontinuity is visible in Figures 18c and 18d (zigzag evolution). Such noise is likely introduced by the way we sampled the contour. Also, the noise in the contour extraction is more visible on Fig. 18g and 18h. The two plots show that our framework can converge to the target in the presence of noisy image processing. When we have false contour data, the value of ASE may encounter a sudden discontinuity. Figure 19 shows examples of these false samples, output by the image processing pipeline. Despite these errors, thanks to the “forgetting nature” of the receding horizon and to the relatively small window size ( $M = 5$ ), the corrupted data will soon be forgotten, and it will not hinder the overall manipulation task. Yet, the overall framework would benefit from a more robust sensing strategy, as in (Chi and Berenson).

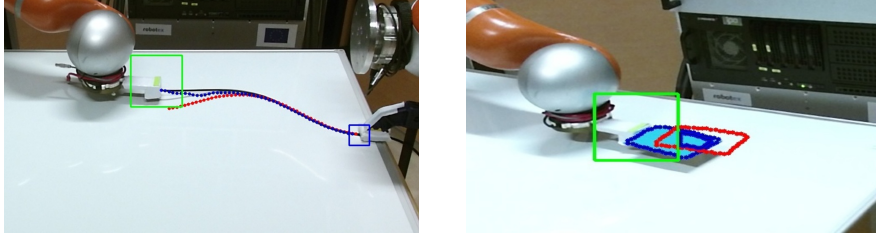


Figure 19: False contour data from the image can cause noise in ASE.

Finally, since our framework can deal with both rigid and deformable objects, we tested it in two experiments where the same object (a sponge) can be both rigid (in the free space), and deformed (when in contact with the environment). These experiments require the robot to: 1) move the object, establish contact, 2) give the object the target shape, by relying on the contact. Figure 20 presents these two original “move and shape” servoing experiments with the corresponding errors ASE plotted in Fig. 21. We use a second fixed robot arm to generate the deforming contact. As the figures and curves show, both experiments were successful.

The success of the “move and shape” task is largely dependent on the contact establishment. However, even when the initial contact has some misalignment (see Fig. 20c and 20g), our framework can still reduce the ASE to give a reasonable final configuration (see Fig. 20h and Fig. 21b).

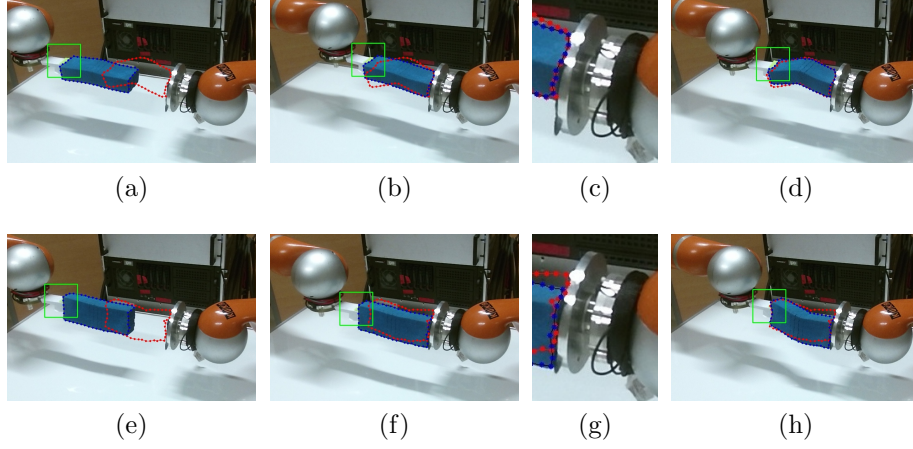


Figure 20: Two “move and shape” experiments grouped into two rows. The target contour (red dotted) is far from the initial one. This requires the robot to 1) move the object, establish contact with the right – fixed – robot arm, 2) give the object the target shape, by relying on the contact. The first column shows the starting configuration, the second column presents the contact establishment, and the third column zooms in to show the alignment. The last column shows the final results.

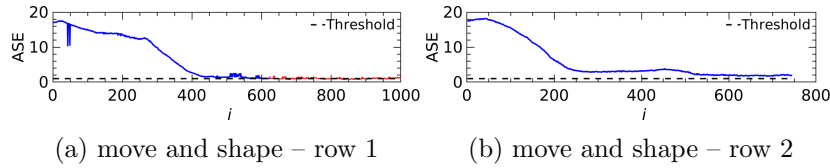


Figure 21: The evolution of  $e_i$  for the experiments of Fig. 20. The black dashed line indicates the threshold  $ASE = 1$ . The blue curves show  $e_i$  until the termination condition, whereas the red curves show the error until manual termination by the human operator.



## 7. Conclusion

2 In this paper, we propose algorithms to automatically and concurrently  
generate object representations (feature vectors) and models of interaction  
4 (interaction matrices) from the same data. We use these algorithms to gen-  
erate the control inputs enabling a robot to move and shape the said object,  
6 be it rigid or deformable. The scheme is validated with comprehensive exper-  
iments, including a target contour that requires both moving and shaping.  
8 We believe it is unprecedented in previous research. Our framework adopts  
a model-free approach. The system characteristics are computed online with  
10 visual and manipulation data. We do not require camera calibration, nor a  
priori knowledge of the camera pose, object size or shape.

12 The proposed approach has two major limitations: 1. *The challenge of*  
*extending it to 6 DoF motion*, 2. *Global convergence cannot be guaranteed*.  
14 Below, we discuss each limitation and present possible solutions.

An open question remains the management of 6 DOF motion of the robot.  
16 Indeed, while the proposed controller can be easily generalized to 6 DOF  
motions, it relies on a sufficiently accurate extraction of feature vectors from  
18 vision sensors. A very challenging task is to generate complete and reliable  
3D feature vectors of objects from a limited sensor set, due to partial views  
20 of the object and to occlusions. To extend it, the framework should benefit  
from robust deformation sensing. In addition, since the approach relies on  
22 local linear models, we expect that with higher DOF, the algorithm will more  
likely get stuck in local minima.

24 The second drawback is that the representation and model of interaction  
are local. Thus, they cannot guarantee global convergence. In addition, our  
26 framework cannot infer whether a shape is reachable or not. This draw-  
back is solvable by using a global deformation model for control. But as we  
28 mentioned earlier, a global model usually requires an offline identification  
phase which we want to avoid. In fact, for different objects, we will need  
30 to re-identify the model. There is a dilemma in using a global deformation  
model.

32 Maybe one of the possible solutions to this dilemma is to have both our  
method and deep learning based methods run in parallel. While our scheme  
34 enables fast online computation and direct manipulation, the extracted data  
can be used by a deep neural network to obtain a global interaction mapping.  
36 Once a global mapping is learned, it can later be used for direct manipulation  
and to infer feasibility of the goal shape.

## Acknowledgement

2 This work is supported in part by the EU H2020 research and innovation  
programme as part of the project VERSATILE under grant agreement No  
4 731330, by the Research Grants Council (RGC) of Hong Kong under grant  
number 14203917, and by the PROCORE-France/Hong Kong RGC Joint  
6 Research Scheme under grant F-PolyU503/18.

## References

- 8 Bakthavatchalam, M., Chaumette, F., Marchand, E., 2013. Photometric  
moments: New promising candidates for visual servoing, in: 2013 IEEE  
10 Int. Conf. on Robotics and Automation, IEEE. pp. 5241–5246.
- Berenson, D., 2013. Manipulation of deformable objects without modeling  
12 and simulating deformation, in: 2013 IEEE/RSJ Int. Conf. on Intelligent  
Robots and Systems, IEEE. pp. 4525–4532.
- 14 Broyden, C.G., 1965. A class of methods for solving nonlinear simultaneous  
equations. *Mathematics of computation* 19(92), 577–593.
- 16 Chaumette, F., 2004. Image moments: a general and useful set of features  
for visual servoing. *IEEE Trans. on Robotics* 20(4), 713–723.
- 18 Chaumette, F., Hutchinson, S., 2006. Visual servo control, part I: Basic  
approaches. *IEEE Robotics and Automation Magazine* 13, 82–90.
- 20 Chaumette, F., Hutchinson, S., 2007. Visual servo control, part II: Advanced  
approaches. *IEEE Robotics and Automation Magazine* 14, 109–118.
- 22 Cherubini, A., Ortenzi, V., Cosgun, A., Lee, R., Corke, P., 2020. Model-free  
vision-based shaping of deformable plastic materials. *The Int. Journal of*  
24 *Robotics Research* 39, 1739–1759.
- Chi, C., Berenson, D., . Occlusion-robust deformable object tracking without  
26 physics simulation, in: 2019 IEEE/RSJ Int. Conf. on Intelligent Robots  
and Systems.
- 28 Collewet, C., Marchand, E., 2011. Photometric visual servoing. *IEEE Trans.*  
*on Robotics* 27, 828–834.

- Collewet, C., Marchand, E., Chaumette, F., 2008. Visual servoing set free  
2 from image processing, in: 2008 IEEE Int. Conf. on Robotics and Automation, IEEE. pp. 81–86.
- 4 Deguchi, K., Noguchi, T., 1996. Visual servoing using eigenspace method  
and dynamic calculation of interaction matrices, in: Proc. of 13th IEEE  
6 Int. Conf. on Pattern Recognition.
- Feller, W., 2008. An introduction to probability theory and its applications.  
8 volume 2. John Wiley & Sons.
- Hosoda, K., Asada, M., 1994. Versatile visual servoing without knowledge of  
10 true jacobian, in: IEEE/RSJ Int. Conf. on Intelligent Robots and Systems.
- Hu, Z., Han, T., Sun, P., Pan, J., Manocha, D., 2019. 3-d deformable object  
12 manipulation using deep neural networks. IEEE Robotics and Automation  
Letters 4, 4255–4261.
- 14 Inoue, H., 1984. Hand-eye coordination in rope handling, in: Proc. of Int.  
Symposium on Robotics Research, MIT PRESS. pp. 163–174.
- 16 Jagersand, M., Fuentes, O., Nelson, R., 1997. Experimental evaluation of  
uncalibrated visual servoing for precision manipulation, in: Int. Conf. on  
18 Robotics and Automation.
- Lagneau, R., Krupa, A., Marchal, M., 2020a. Active deformation through  
20 visual servoing of soft objects, in: IEEE Int. Conf. on Robotics and Automation.
- 22 Lagneau, R., Krupa, A., Marchal, M., 2020b. Automatic shape control of  
deformable wires based on model-free visual servoing. IEEE Robotics and  
24 Automation Letters 5, 5252–5259.
- Lapresté, J.T., Jurie, F., Dhome, M., Chaumette, F., 2004. An efficient  
26 method to compute the inverse jacobian matrix in visual servoing, in:  
IEEE Int. Conf. on Robotics and Automation.
- 28 Laranjeira, M., Dune, C., Hugel, V., 2017. Catenary-based visual servoing  
for tethered robots, in: 2017 IEEE Int. Conf. on Robotics and Automation  
30 (ICRA), IEEE. pp. 732–738.

- Laranjeira, M., Dune, C., Hugel, V., 2020. Catenary-based visual servoing for tether shape control between underwater vehicles. *Ocean Engineering* 200, 1–19.
- Li, X., Su, X., Gao, Y., Liu, Y.H., 2018. Vision-based robotic grasping and manipulation of usb wires, in: *IEEE Int. Conf. on Robotics and Automation (ICRA)*.
- Marchand, E., 2019. Subspace-based direct visual servoing. *IEEE Robotics and Automation Letters* 4, 2699–2706.
- Nair, A., Chen, D., Agrawal, P., Isola, P., Abbeel, P., Malik, J., Levine, S., 2017. Combining self-supervised learning and imitation for vision-based rope manipulation, in: *IEEE Int. Conf. on Robotics and Automation*.
- Navarro-Alarcon, D., Cherubini, A., Li, X., 2019. On model adaptation for sensorimotor control of robots, in: *Chinese Control Conference*, pp. 2548–2552.
- Navarro-Alarcon, D., Liu, Y., 2013. Uncalibrated vision-based deformation control of compliant objects with online estimation of the jacobian matrix, in: *IEEE/RSJ Int. Conf. on Intelligent Robots and Systems*.
- Navarro-Alarcon, D., Liu, Y., Romero, J.G., Li, P., 2013. Model-free visually servoed deformation control of elastic objects by robot manipulators. *IEEE Trans. on Robotics* 29(6), 1457–1468.
- Navarro-Alarcon, D., Liu, Y., Romero, J.G., Li, P., 2014. On the visual deformation servoing of compliant objects: Uncalibrated control methods and experiments. *Int. Journal of Robotics Research* 33(11), 1462–1480.
- Navarro-Alarcon, D., Liu, Y.H., 2018. Fourier-based shape servoing: A new feedback method to actively deform soft objects into desired 2d image contours. *IEEE Trans. on Robotics* 34(1), 272–279.
- Nayar, S.K., Nene, S.A., Murase, H., 1996. Subspace methods for robot vision. *IEEE Trans. on Robotics and Automation* 12(5), 750–758.
- Sanchez, J., Corrales, J.A., Bouzgarrou, B.C., Mezouar, Y., 2018. Robotic manipulation and sensing of deformable objects in domestic and industrial applications: a survey. *The Int. Journal of Robotics Research* .

- 2 Sang, Q., Tao, G., 2012. Adaptive control of piecewise linear systems: the  
state tracking case. *IEEE Trans. Autom. Control* 57, 522–528.
- 4 She, Y., Wang, S., Dong, S., Sunil, N., Rodriguez, A., Adelson, E., 2020.  
Cable manipulation with a tactile-reactive gripper, in: *Robotics: Science  
and Systems*.
- 6 Smith, P.W., Nandhakumar, N., Ramadorai, A.K., 1996. Vision based ma-  
nipulation of non-rigid objects, in: *IEEE Int. Conf. on Robotics and Au-  
8 tomation*, IEEE. pp. 3191–3196.
- 10 Wakamatsu, H., Hirai, S., 2004. Static modeling of linear object deformation  
based on differential geometry. *The Int. Journal of Robotics Research* 23,  
293–311.
- 12 Wakamatsu, H., Hirai, S., Iwata, K., 1995. Modeling of linear objects con-  
sidering bend, twist, and extensional deformations, in: *IEEE Int. Conf. on  
14 Robotics and Automation*.
- 16 Wang, Z., Li, X., Navarro-Alarcon, D., Liu, Y.h., . A unified controller for  
region-reaching and deforming of soft objects, in: *2018 IEEE/RSJ Int.  
Conf. on Intelligent Robots and Systems (IROS)*.
- 18 Zeng, R., Wu, J., Shao, Z., Chen, Y., Chen, B., Senhadji, L., Shu, H., 2016.  
Color image classification via quaternion principal component analysis net-  
20 work. *Neurocomputing* 216, 416–428.
- 22 Zhang, L., Dong, W., Zhang, D., Shi, G., 2010. Two-stage image denoising by  
principal component analysis with local pixel grouping. *Pattern recognition*  
43, 1531–1549.
- 24 Zhu, J., Navarro, B., Fraitse, P., Crosnier, A., Cherubini, A., 2018. Dual-  
arm robotic manipulation of flexible cables, in: *IEEE/RSJ Int. Conf. on  
26 Intelligent Robots and Systems*.