

Model Predictive Manipulation of Compliant Objects with Multi-Objective Optimizer and Adversarial Network for Occlusion Compensation

Jiaming Qi, Dongyu Li, Yufeng Gao, Peng Zhou, and David Navarro-Alarcon, *Senior Member, IEEE*

Abstract—The robotic manipulation of compliant objects is currently one of the most active problems in robotics due to its potential to automate many important applications. Despite the progress achieved by the robotics community in recent years, the 3D shaping of these types of materials remains an open research problem. In this paper, we propose a new vision-based controller to automatically regulate the shape of compliant objects with robotic arms. Our method uses an efficient online surface/curve fitting algorithm that quantifies the object’s geometry with a compact vector of features; This feedback-like vector enables to establish an explicit shape servo-loop. To coordinate the motion of the robot with the computed shape features, we propose a receding-time estimator that approximates the system’s sensorimotor model while satisfying various performance criteria. A deep adversarial network is developed to robustly compensate for visual occlusions in the camera’s field of view, which enables to guide the shaping task even with partial observations of the object. Model predictive control is utilized to compute the robot’s shaping motions subject to workspace and saturation constraints. A detailed experimental study is presented to validate the effectiveness of the proposed control framework.

Note to Practitioners—This paper is motivated by the problem of manipulating compliant objects by robotic arms. Classical vision-based controllers are unsuitable for this task as they rely on ideal setups that are difficult to meet in practice, e.g., unoccluded views and unconstrained environments. The method proposed in this paper aims to address these limitations. For that, we use curve fitting algorithms to compute features that approximate the object’s geometry; This approach is more efficient than traditional dense point clouds. To facilitate the method’s implementation in real-world environments, we propose an optimization-based controller that computes the robot’s motion while simultaneously satisfying various constraints, e.g., workspace limits, input saturation, etc. The reported experiments validate the feasibility of our strategy in dynamic/unstructured situations, yet, note that this approach may not properly shape materials with negligible elastic properties (e.g., fabrics, food materials). Our proposed methods can be used in many other applications ranging from classical pick-and-place tasks to visually guiding robots with unknown kinematic models.

Index Terms—Robotics; Visual Servoing; Deformable Objects; Occlusion Compensation; Model Predictive Control

This work is supported in part by the Research Grants Council (RGC) of Hong Kong under grants 14203917 and 15212721, in part by the Key-Area Research and Development Program of Guangdong Province under Grant 2020B090928001, and in part by the Jiangsu Industrial Technology Research Institute Collaborative Funding Scheme under grant 43-ZG9V.

J. Qi and Y. Gao are with the Harbin Institute of Technology, Department of Control Science and Engineering, Harbin, Heilongjiang, China. (e-mail: 18B904030@hit.edu.cn, gaoyf@stu.hit.edu.cn)

D. Li is with Beihang University, School of Cyber Science and Technology, Beijing, China. (e-mail: dongyuli@buaa.edu.cn)

P. Zhou and D. Navarro-Alarcon are with The Hong Kong Polytechnic University, Department of Mechanical Engineering, Kowloon, Hong Kong. (e-mail: jeffery.zhou@connect.polyu.hk, dna@ieee.org)

I. INTRODUCTION

THE manipulation/shaping of deformable bodies by robots is a fundamental problem that has recently attracted the attention of many researchers [1]; Its complexity has forced researchers to develop new methods in a wide range of fundamental areas that include representation, learning, planning, and control. From an applied research perspective, this challenging problem has shown great potential in various economically-important tasks such as the assembly of compliant/delicate objects [2], surgical/medical robotics [3], cloth/fabric folding [4], etc. The manipulation of compliant materials contrast with its rigid-body counterpart in that physical interactions will invariably change the object’s shape, which introduces additional degrees-of-freedom to the typically unknown objects, and hence, complicates the manipulation task. While great progress has been achieved in recent years, the development of these types of embodied manipulation capabilities is still largely considered an open research problem in robotics and control.

There are various technical issues that hamper the implementation of these tasks in real-world unstructured environments, which largely differ from implementations in ideal simulation environments (a trend followed by many recent works). Here, we argue that to effectively servo-control the shape of deformable materials in the field, a sensor-based controller must possess the following features: 1) Efficient compression of the object’s high-dimensional shape; 2) Occlusion-tolerant estimation of the objects geometry; 3) Adaptive prediction of the differential shape-motion model; Our aim in this paper is precisely to develop a new shape controller endowed with all the above-mentioned features. The proposed method is formulated under the model predictive control (MPC) framework that enables to compute shaping actions that satisfy multiple performance criteria (a key property for real engineering applications).

A. Related Work

Many researchers have previously studied this challenging problem (we refer the reader to [5], [6] for comprehensive reviews). To servo-control the object’s non-rigid shape, it is essential to design a low-dimensional feature representation that can capture the key geometric properties of the object. Several representation methods have been proposed before, e.g., geometric features based on points, angles, curvatures, etc [7], [8]; However, due to their hard-coded nature, these methods can only be used to represent a single shaping

action. Other geometric features computed from contours and centerlines [9]–[11] can represent soft object deformation in a more general way. Various data-driven approaches have also been proposed to represent shapes, e.g., using fast point feature histograms [12], bottleneck layers [13], [14], principal component analysis [15], etc. However, there is no widely accepted approach to compute efficient/compact feature representations for 3D shape; This is still an open research problem.

Occlusions of a camera’s field of view pose many complications to the implementation of visual servoing controllers, as the computation of (standard) feedback features requires complete visual observations of the object at all times. Many methods have been developed to tackle this critical issue, e.g. [16] used the estimated interaction matrix to handle information loss in visual servoing, yet, this method requires a calibrated visual-motor model. Coherent point drift was utilized in [17] to register the topological structure from previous sequences and to predict occlusions; However, this method is sensitive to the initial point set, further affects the registration result. A structure preserved registration method was presented in [18] to track occluded objects; This approach has good accuracy and robustness to noise, yet, its efficiency decreases with the number of points. To efficiently implement vision-based strategies in the field, it is essential to develop algorithms that can robustly guide the servoing task, even in the presence of occlusions.

To visually guide the manipulation task, control methods must have some form of model that (at least approximately) describes that how the input robot motions produce output shape changes [19]; In the visual servoing community, such differential relation is typically captured by the so-called interaction (Jacobian) matrix [20]. Many methods have been proposed to address this issue, e.g. the Broyden update rule [21] is a classical algorithm to iteratively estimate this transformation matrix [7], [22], [23]. Although these types of algorithms do not require knowledge of the model’s structure, its estimation properties are only valid locally. Other approaches with global estimation properties include algorithms based on (deep) artificial neural networks [2], optimization-based algorithms [24], adaptive estimators [25], etc. However, the majority of existing methods estimate this model based on a single performance criterion (typically, a first-order Jacobian-like relation), which limits the types of dynamic responses that the robot can achieve during a task. A multi-objective model estimation is particularly important in the manipulation of compliant materials, as their mechanical properties are rarely known in practice, thus, making hard to meet various performance requirements.

To compute the active shaping motions, most control methods only formulate the problem in terms of the final target shape and do not typically consider the system’s physical constraints. MPC represents a feasible solution to these issues, as it performs the control tasks by optimizing cost functions over a finite time-horizon rather than finding the exact analytical solution [26]; This allows MPC to compute controls that guide the task while satisfying a set of constraints, e.g., control saturation, workspace bounds, etc. Despite its valuable and flexible properties, MPC has not been sufficiently studied in the context of shape control of deformable objects.

B. Our Contribution

To solve the above-mentioned issues, in this paper we propose a new control framework to manipulate purely-elastic objects into desired configurations. The original features of our new methodology include: 1) A parametric shape descriptor to efficiently characterize 3D deformations based on online curve/surface fitting; 2) A robust shape prediction network based on adversarial neural networks to compensate visual occlusions; 3) An optimization-based estimator to approximate the deformation Jacobian matrix and satisfy various performance constraints; 4) An MPC-based motion controller to guide the shaping motions while simultaneously solving workspace and saturation constraints.

To the best of the authors’ knowledge, this is the first time that a shape servoing controller is developed with all the functions proposed in this paper. To validate the effectiveness of our new methodology, we report a detail experimental study with a robotic platform manipulating various types of compliant objects.

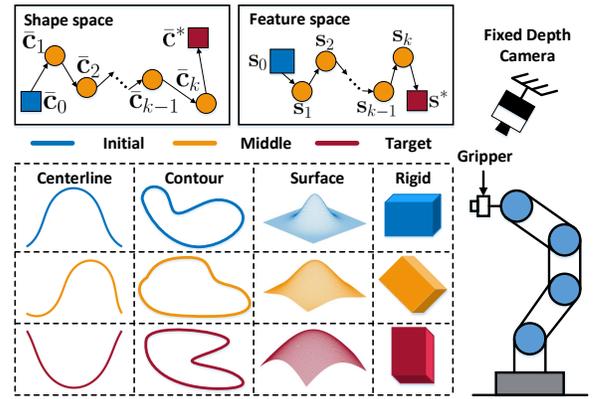


Fig. 1. Schematic diagram of the manipulations of elastic objects, including deformation (centerline, contour, and surface) and positioning tasks of rigid objects. The framework aims to command the robot to manipulate elastic objects into the target configurations.

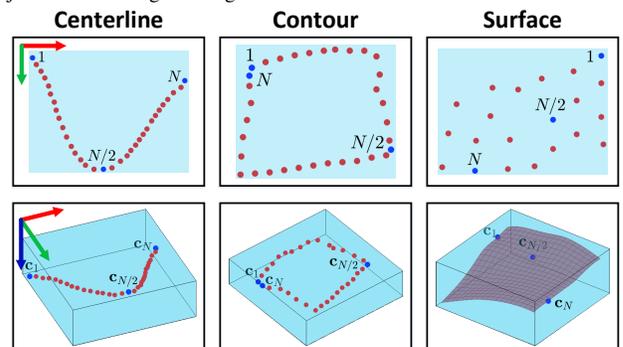


Fig. 2. Various shape configurations, including centerline, contour and surface. The first row shows the pixel coordinates for each shape. The second row shows the Cartesian coordinates related to the corresponding pixel ones through *OpenCV/RealSense*. The generated shapes are ordered, fixed-sampled and equidistant. The rigid object adopts surface representation.

II. PROBLEM FORMULATION

Notation: In this paper, we use the following frequently-used notation: Bold small letters, e.g., \mathbf{v} , denote column vectors, while bold capital letters, e.g., \mathbf{M} , denote matrices. Time evolving variables are denoted as \mathbf{x}_k , for k as the discrete time instant. The $n \times m$ matrix of ones is denoted by $\mathbf{I}_{n \times m}$

and the identity matrix as \mathbf{E}_n . \mathbf{L}_n represents the low triangle matrix of $\mathbf{I}_{n \times n}$, and \otimes represents the Kronecker product.

The schematic diagram of the proposed shape servoing framework is conceptually illustrated in Fig. 1. A depth camera with eye-to-hand configuration observes the shapes of elastic objects that are manipulated by the robot, see Fig. 2. We denote the 3D measurement points captured by the vision system as:

$$\bar{\mathbf{c}} = [\mathbf{c}_1^\top, \dots, \mathbf{c}_N^\top]^\top \in \mathbb{R}^{3N}, \quad \mathbf{c}_i = [x_i, y_i, z_i]^\top \in \mathbb{R}^3 \quad (1)$$

for N as the number of points, and \mathbf{c}_i as the 3D coordinates of the i th point, expressed in the camera frame.

A. Feature Sensorimotor Model

Let us denote the position of the robot's end-effector by $\mathbf{r} \in \mathbb{R}^3$. As the dimension of the $3N$ observed points $\bar{\mathbf{c}}$ is typically very large, therefore, its direct use as a feedback signal for shape servocontrol is impractical. Therefore, an efficient controller may typically require to use some form of dimension reduction technique. To deal with this issue, we construct a feature vector $\mathbf{s} = \mathbf{f}_s(\bar{\mathbf{c}}) : \mathbb{R}^{3N} \mapsto \mathbb{R}^p$, for $p \ll 3N$, to represent the geometric feedback $\bar{\mathbf{c}}$ of the object. This compact feedback-like signal will be used to design our automatic 3D shape controller.

For *purely elastic* objects undergoing deformations, it is reasonable to model that the object configuration is only dependant on its potential energy \mathcal{P} (thus, all inertial and viscous effects are neglected from our analysis [10]). We model that \mathcal{P} is fully determined by the feedback feature vector \mathbf{s} and the robot's position \mathbf{r} [27], i.e.: $\mathcal{P} = \mathcal{P}(\mathbf{s}, \mathbf{r})$. In steady-state, the extremum expression satisfies $(\partial \mathcal{P} / \partial \mathbf{s})^\top = \mathbf{a}(\mathbf{s}, \mathbf{r}) = \mathbf{0}$ [28]. We then define the following matrices:

$$\frac{\partial^2 \mathcal{P}}{\partial \mathbf{s}^2} = \mathbf{G}(\mathbf{s}, \mathbf{r}), \quad \frac{\partial^2 \mathcal{P}}{\partial \mathbf{r} \partial \mathbf{s}} = \mathbf{K}(\mathbf{s}, \mathbf{r}) \quad (2)$$

which are useful to linearize the extremum equation (by using first-order Taylor's series expansion) as follows:

$$\mathbf{a}(\mathbf{s} + \Delta \mathbf{s}, \mathbf{r} + \Delta \mathbf{r}) \approx \mathbf{a}(\mathbf{s}, \mathbf{r}) + \mathbf{G}(\mathbf{s}, \mathbf{r}) \Delta \mathbf{s} + \mathbf{K}(\mathbf{s}, \mathbf{r}) \Delta \mathbf{r} \quad (3)$$

for $\Delta \mathbf{s}$ and $\Delta \mathbf{r}$ as small changes. Note that as $\mathbf{a}(\mathbf{s} + \Delta \mathbf{s}, \mathbf{r} + \Delta \mathbf{r}) = \mathbf{a}(\mathbf{s}, \mathbf{r}) = \mathbf{0}$ is satisfied, we can obtain the following motion model:

$$\Delta \mathbf{s} \approx -\mathbf{G}(\mathbf{s}, \mathbf{r})^{-1} \mathbf{K}(\mathbf{s}, \mathbf{r}) \Delta \mathbf{r} = \mathbf{J}(\mathbf{s}, \mathbf{r}) \mathbf{u} \quad (4)$$

where $\mathbf{J}(\mathbf{s}, \mathbf{r}) = -\mathbf{G}(\mathbf{s}, \mathbf{r})^{-1} \mathbf{K}(\mathbf{s}, \mathbf{r}) \in \mathbb{R}^{p \times 3}$ represents the deformation Jacobian matrix (which depends on both the feature vector and robot position), and \mathbf{u} represents the robot's motion control input. This model can be expressed in an intuitive discrete-time form:

$$\Delta \mathbf{s}_{k+1} = \mathbf{J}_k(\mathbf{s}_k, \mathbf{r}_k) \mathbf{u}_k \quad (5)$$

The deformation Jacobian matrix (DJM) \mathbf{J}_k indicates how the robot's action $\mathbf{u}_k = \mathbf{r}_k - \mathbf{r}_{k-1}$ produces changes in the feedback features $\Delta \mathbf{s}_{k+1} = \mathbf{s}_{k+1} - \mathbf{s}_k$. Clearly, the analytical computation of \mathbf{J}_k requires knowledge of the physical properties and model of the elastic object and the vision system, which are difficult to obtain in practice. Thus, numerical methods are often used to approximate this matrix in real-time, which enables to perform vision-guided manipulation tasks.

In this paper, we consider a robot manipulator whose control inputs represent velocity commands (here, modelled as the differential changes $\Delta \mathbf{r}$). It is assumed that $\Delta \mathbf{r}$ can be instantaneously executed without delay [11].

Problem statement. Design a vision-based control method to automatically manipulate a compliant object into a target 3D configuration, while simultaneously compensating for visual occlusions of the camera and estimating the deformation Jacobian matrix of the object-robot system.

III. SHAPE REPRESENTATION

This section presents how online curve/surface fitting (least-squares minimization (LSM) [29] and moving least squares (MLS) [30]) is combined with a parametric shape descriptor to compute a compact vector of feedback shape features.

A. LSM-Based Features

1) *Centerline and Contour Extraction:* Both are expressed as a parametric curve dependent on the normalized arc-length $0 \leq \rho \leq 1$. Then, the point can be represented as $\mathbf{c}_i = \mathbf{f}(\rho_i)$, with ρ_i as the arc-length between the start point \mathbf{c}_1 and \mathbf{c}_i , where $\rho_1 = 0$ and $\rho_N = 1$. The fitting functional $\mathbf{f}(\cdot)$ is constructed as follows:

$$\mathbf{f}(\rho) = \sum_{j=0}^n \mathbf{p}_j B_{j,n}(\rho) \quad (6)$$

where the vector $\mathbf{p}_j \in \mathbb{R}^3$ denotes the shape weights, $n \in \mathbb{N}^*$ specifies the fitting order, and the scalar $B_{j,n}(\rho)$ represents a parametric regression function, which may take various forms, such as: [31]

- Polynomial parameterization [32]:

$$B_{j,n}(\rho) = \rho^j \quad (7)$$

- Bernstein parameterization [33]:

$$B_{j,n}(\rho) = C_n^j (1 - \rho)^{n-j} \rho^j \quad (8)$$

where C_b^j represents the binomial coefficient.

- Cox-deBoor parameterization [34]:

$$B_{j,n}(\rho) = \frac{1}{n!} \sum_{l=0}^{n-j} \left((-1)^l C_{n+1}^l (\rho + n - j - l)^n \right) \quad (9)$$

- Trigonometric parameterization [35]:

$$B_{j,n}(\rho) = \begin{cases} 1, & j = 0 \\ \cos(\frac{j+1}{2}\rho), & j > 0, j \text{ is odd} \\ \sin(\frac{j}{2}\rho), & j > 0, j \text{ is even} \end{cases} \quad (10)$$

From (1) and (6), we can compute the following fitting cost function:

$$Q = (\mathbf{B}\mathbf{s} - \bar{\mathbf{c}})^\top (\mathbf{B}\mathbf{s} - \bar{\mathbf{c}}) \quad (11)$$

for a "tall" regression-like matrix \mathbf{B} constructed as:

$$\mathbf{B} = [\mathbf{B}_1^\top, \dots, \mathbf{B}_N^\top]^\top \in \mathbb{R}^{3N \times 3(n+1)} \\ \mathbf{B}_i = [B_{0,n}(\rho_i), \dots, B_{n,n}(\rho_i)] \otimes \mathbf{E}_3 \in \mathbb{R}^{3 \times 3(n+1)} \quad (12)$$

and $\mathbf{s} = [\mathbf{p}_0^\top, \dots, \mathbf{p}_n^\top]^\top \in \mathbb{R}^{3(n+1)}$ as a compact feedback feature vector that represents the shape. We seek to minimize

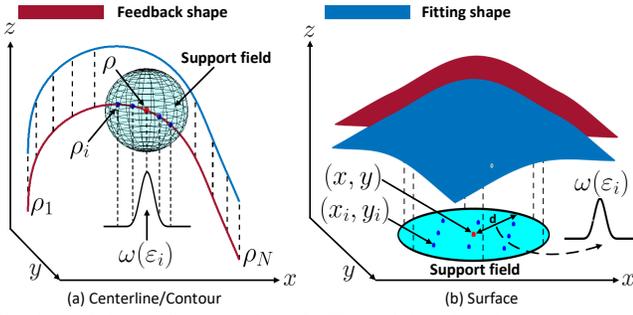


Fig. 3. Scheme diagram of MLS. The weight $\omega(\varepsilon_i)$ is the square error between the fitted value and the given value. The fitting smoothness is regulated by adjusting the weight values.

(11) to obtain a feature vector \mathbf{s} that closely approximates $\bar{\mathbf{c}} \approx \mathbf{B}\mathbf{s}$. The solution to (11) is:

$$\mathbf{s} = (\mathbf{B}^T \mathbf{B})^{-1} \mathbf{B}^T \bar{\mathbf{c}} \quad (13)$$

where it is assumed that $N \gg n + 1$.

2) *Surface Extraction*: The equation of the surface is defined as follows:

$$z = f(x, y) = \sum_{j=0}^{n_x} \sum_{l=0}^{n_y} B_{j,n_x}(x) B_{l,n_y}(y) q_{jl} \quad (14)$$

where $n_x, n_y \in \mathbb{N}^*$ are the fitting order along with x and y direction, and $q_{jl} \in \mathbb{R}$ is the shape weight. Same as with (11), as fitting cost function is introduced:

$$Q = (\mathbf{D}\mathbf{s} - \mathbf{z})^T (\mathbf{D}\mathbf{s} - \mathbf{z}) \quad (15)$$

for a ‘‘augmented’’ regression matrix \mathbf{D} satisfying:

$$\begin{aligned} \mathbf{D} &= [\mathbf{D}_1^T, \dots, \mathbf{D}_N^T]^T \in \mathbb{R}^{N \times (n_x+1)(n_y+1)} \\ \mathbf{D}_i &= [B_{0,n_x}(x_i) B_{0,n_y}(y_i), \dots, B_{n_x,n_x}(x_i) B_{n_y,n_y}(y_i)] \\ \mathbf{D}_i^T &\in \mathbb{R}^{(n_x+1)(n_y+1)}, \text{ for } i = 1, \dots, N \end{aligned} \quad (16)$$

The depth vector is defined as $\mathbf{z} = [z_1, \dots, z_N]^T \in \mathbb{R}^N$, and the feature vector is $\mathbf{s} = [q_{00}, \dots, q_{n_x n_y}]^T \in \mathbb{R}^{(n_x+1)(n_y+1)}$. The solution to the minimization of (15) that approximates $\mathbf{z} \approx \mathbf{D}\mathbf{s}$ is as follows:

$$\mathbf{s} = (\mathbf{D}^T \mathbf{D})^{-1} \mathbf{D}^T \mathbf{z} \quad (17)$$

For $N \gg (n_x + 1)(n_y + 1)$.

B. MLS-based Feature Extraction

Although LSM has an efficient one-step calculation, the weights of $B_{j,n}$ are the same all over the variable parameters (e.g., ρ or x, y). Thus, to approximate complex curves/surfaces, the fitting order needs to be increased, which may lead to overfitting problems. To address this issue, MLS assumes that \mathbf{s} is parameter-dependent, i.e., it changes with respect to ρ (for centerline and contour), or to (x, y) (for surface). This enables MLS to represent complex shapes with a lower fitting order. A support field [36] is introduced to ensure that the value of each weight is only influenced by data points within the support field. A conceptual diagram is given in Fig. 3.

¹In the following sections, Q is used to generically represent different, albeit related, functions.

1) *Centerline and Contour Extraction*: The parametric equation with $\sigma_j(\rho)$ is presented by referring to (6):

$$\mathbf{f}(\rho) = \sum_{j=0}^n \sigma_j(\rho) B_{j,n}(\rho) \quad (18)$$

where the parameter-dependent weight is denoted as $\sigma_j(\rho) \in \mathbb{R}^3$. The weighted square residual function is defined as:

$$Q(\rho) = \sum_{i=1}^N \omega(\varepsilon_i) \left\| \sum_{j=0}^n \sigma_j(\rho) B_{j,n}(\rho_i) - \mathbf{c}_i \right\|^2 \quad (19)$$

where $\varepsilon_i = |\rho - \rho_i|/d > 0$, and d is the constant support field radius. The scalar function $\omega(\varepsilon_i)$ is calculated as follows [37]:

$$\omega(\varepsilon_i) = \begin{cases} \frac{2}{3} - 4\varepsilon_i^2 + 4\varepsilon_i^3, & \varepsilon_i \leq 0.5 \\ \frac{4}{3} - 4\varepsilon_i + 4\varepsilon_i^2 - \frac{4}{3}\varepsilon_i^3, & 0.5 < \varepsilon_i \leq 1 \\ 0, & \varepsilon_i > 1 \end{cases} \quad (20)$$

The function $\omega(\varepsilon_i)$ indicates the weight of ρ relative to ρ_i , $\omega(\varepsilon_i)$ decreases as ε_i increasing inside the support field. When ρ is outside the support field, $\omega(\varepsilon_i) = 0$. MLS reduces into LSM when $\omega(\varepsilon_i)$ is constant. The cost (19) can be equivalently expressed in matrix form as:

$$Q(\rho) = (\mathbf{B}\boldsymbol{\vartheta}(\rho) - \bar{\mathbf{c}})^T \mathbf{W}(\varepsilon) (\mathbf{B}\boldsymbol{\vartheta}(\rho) - \bar{\mathbf{c}}) \quad (21)$$

where $\boldsymbol{\vartheta}(\rho)$ and $\mathbf{W}(\varepsilon)$ are defined as follows:

$$\begin{aligned} \boldsymbol{\vartheta}(\rho) &= [\sigma_0^T(\rho), \dots, \sigma_n^T(\rho)]^T \in \mathbb{R}^{3(n+1)} \\ \mathbf{W}(\varepsilon) &= \text{diag}(\omega(\varepsilon_1), \dots, \omega(\varepsilon_N)) \otimes \mathbf{E}_3 \in \mathbb{R}^{3N \times 3N} \end{aligned} \quad (22)$$

The value of $\boldsymbol{\vartheta}(\rho)$ that minimizes (21) is computed as follows:

$$\boldsymbol{\vartheta}(\rho) = (\mathbf{B}^T \mathbf{W}(\varepsilon) \mathbf{B})^{-1} \mathbf{B}^T \mathbf{W}(\varepsilon) \bar{\mathbf{c}} \quad (23)$$

By completing N iterations along the parameter ρ_1, \dots, ρ_N , we can compute the augmented shape features as:

$$\boldsymbol{\Pi} = [\boldsymbol{\vartheta}(\rho_1), \dots, \boldsymbol{\vartheta}(\rho_N)] \in \mathbb{R}^{3(n+1) \times N} \quad (24)$$

As the dimension of (24) is very large, it is impractical to use its components as a feedback signal for control. Thus, principal components analysis (PCA) [15] is used to reduce the augmented structure (24) into a compact form $\tilde{\boldsymbol{\Pi}} \in \mathbb{R}^{3(n+1) \times m}$ where $m \ll N$ by selecting the first m most significant dimensions. The feature vector $\mathbf{s} \in \mathbb{R}^{3m(n+1)}$ is computed by vectorizing the elements of $\tilde{\boldsymbol{\Pi}}$.

2) *Surface Extraction*: The equation of the surface is constructed as:

$$f(x, y) = \sum_{j=0}^{n_x} \sum_{l=0}^{n_y} B_{j,n_x}(x) B_{l,n_y}(y) \varpi_{jl}(x, y) \quad (25)$$

where $\varpi_{jl}(x, y) \in \mathbb{R}$ is the parameter-dependent weight related to (x, y) . Algorithm 1 gives a pseudocode description of this method.

Remark 1. In addition to the proposed basis functions, there are other approaches that can be used to obtain similar results, e.g., Chebyshev polynomials [38] and Legendre basis transformations [39].

Algorithm 1 Surface fitting procedure of MLS.

Require: $\bar{\mathbf{c}}, n_x, n_y, m$, and d ;

1: Calculate node distance:

$$\varepsilon_i = \sqrt{(x - x_i)^2 + (y - y_i)^2} / d$$

2: Construct the fitting cost function:

$$Q = (\mathbf{D}\phi(x, y) - \mathbf{z})^T \mathbf{W}(\varepsilon) (\mathbf{D}\phi(x, y) - \mathbf{z})$$

 for $\phi(x, y)$ and $\mathbf{W}(\varepsilon)$ satisfying:

$$\begin{aligned} \phi &= [\varpi_{00}, \dots, \varpi_{n_x n_y}]^T \\ \mathbf{W} &= \text{diag}(\omega(\varepsilon_1), \dots, \omega(\varepsilon_N)) \end{aligned} \quad (26)$$

 3: Compute the structure $\phi(x, y)$:

$$\phi = (\mathbf{D}^T \mathbf{W}(\varepsilon) \mathbf{D})^{-1} \mathbf{D}^T \mathbf{W}(\varepsilon) \mathbf{z}$$

 4: Use $x_i, y_i, i \in [1, N]$ to compute:

$$\mathbf{\Pi} = [\phi(x_1, y_1), \dots, \phi(x_N, y_N)] \quad (27)$$

 5: Use PCA to calculate $\tilde{\mathbf{\Pi}}$;

 6: Vectorize $\tilde{\mathbf{\Pi}}$ to obtain \mathbf{s} ;

 7: **return** \mathbf{s} ;

IV. SHAPE PREDICTION NETWORK

During the manipulation process, occlusions caused by obstacles or the robot itself may affect the integrity of observed shapes, and hence, the vector \mathbf{s} cannot properly describe the object's configuration. As a solution to this critical issue, in this paper we propose an occlusion compensation shape prediction network (SPN), which is composed of a regulation input (RI), a multi-resolution encoder (MRE) and a discriminator network (DN) [40]. The proposed SPN utilizes the robot and object configurations and the active robot motions (i.e., $\bar{\mathbf{c}}_k, \mathbf{r}_k, \mathbf{u}_k$) as input to the network to predict the next instance shape, here denoted by $\hat{\mathbf{c}}_{k+1}$. Fig. 4 shows the overall architecture of the SPN.

A. Input Data Preprocessing

As the input data $\bar{\mathbf{c}}_k, \mathbf{r}_k, \mathbf{u}_k$ to the network have different sizes, they need to be rearranged into structures with unified dimensions. To this end, $\bar{\mathbf{c}}_k$ is first rearranged into the matrix $\mathbf{T}_k^{high} = [\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_N]^T \in \mathbb{R}^{N \times 3}$. Then, farthest point sampling (FPS) [41] is used to downsample \mathbf{T}_k^{high} to two resolutions $\mathbf{T}_k^{mid} \in \mathbb{R}^{\frac{N}{\delta} \times 3}$ and $\mathbf{T}_k^{low} \in \mathbb{R}^{\frac{N}{\delta^2} \times 3}$, for δ as the resolution scale. Finally, the vectors \mathbf{r}_k and \mathbf{u}_k are rearranged into the matrices $\mathbf{\Lambda}_k^{high} = \mathbf{I}_{N \times 1} \otimes \mathbf{r}_k^T \in \mathbb{R}^{N \times 3}$ and $\mathbf{\Sigma}_k^{high} = \mathbf{I}_{N \times 1} \otimes \mathbf{u}_k^T \in \mathbb{R}^{N \times 3}$, which are similarly downsampled into mid and low resolutions as follows:

$$\mathbf{\Lambda}_k^{mid} = \mathbf{I}_{\frac{N}{\delta} \times 1} \otimes \mathbf{r}_k^T, \quad \mathbf{\Lambda}_k^{low} = \mathbf{I}_{\frac{N}{\delta^2} \times 1} \otimes \mathbf{r}_k^T, \quad (28)$$

$$\mathbf{\Sigma}_k^{mid} = \mathbf{I}_{\frac{N}{\delta} \times 1} \otimes \mathbf{u}_k^T, \quad \mathbf{\Sigma}_k^{low} = \mathbf{I}_{\frac{N}{\delta^2} \times 1} \otimes \mathbf{u}_k^T \quad (29)$$

Thus, three different resolutions are generated for the network, high $\{\mathbf{T}_k^{high}, \mathbf{\Lambda}_k^{high}, \mathbf{\Sigma}_k^{high}\}$, mid $\{\mathbf{T}_k^{mid}, \mathbf{\Lambda}_k^{mid}, \mathbf{\Sigma}_k^{mid}\}$, and low $\{\mathbf{T}_k^{low}, \mathbf{\Lambda}_k^{low}, \mathbf{\Sigma}_k^{low}\}$, that provides a total input data of dimension $[N \times 9, \frac{N}{\delta} \times 9, \frac{N}{\delta^2} \times 9]$. $\mathbf{T}_k^{high}, \mathbf{T}_k^{mid}, \mathbf{T}_k^{low}$ are the geometric shapes of the object under different compression

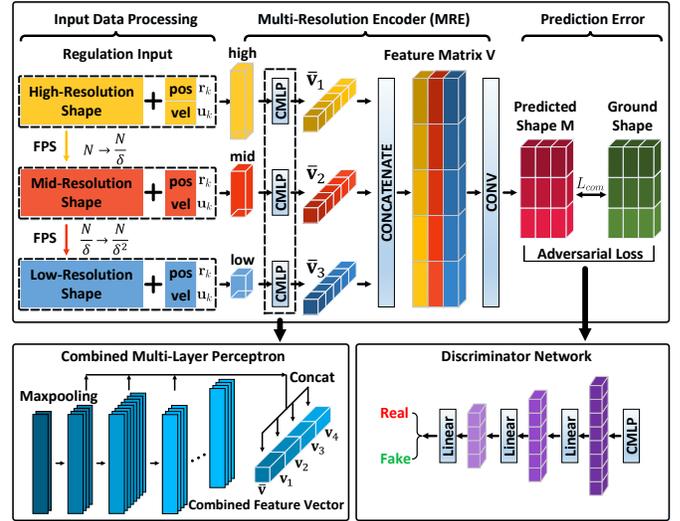


Fig. 4. SPN predicts the next-moment shape $\hat{\mathbf{c}}_{k+1}$ by using the current-moment data in the case of occlusion, i.e., $\bar{\mathbf{c}}_k + \mathbf{r}_k + \mathbf{u}_k \rightarrow \hat{\mathbf{c}}_{k+1}$. FPS [41] regulates $\bar{\mathbf{c}}_k$ (shown in yellow) uniformly to different scales (shown in red and blue). MRE stacks three-resolutions data together to form the total feature information, and obtains the predicted next-moment shape through convolution. DN is used to further improve the accuracy of the network.

sizes specified by δ . Thus, these three terms can describe the potential feature structure of objects. The proposed SPN aims to predict the object's shape that results from the robots actions, under the current object-robot configuration. By using this multi-resolution data $\{\text{high, mid, low}\}$, the encoder can better learn the potential feature information of shapes.

B. Multi-Resolution Encoder

A combined multi-layer perceptron (CMLP) is the feature extractor of MRE, which uses the output of each layer in a MLP to form a multiple-dimensional feature vector. Traditional methods adopt the last layer of the MLP output as features, and do not consider the output of the intermediate layers, which leads to potentially losing important local information [42]. CMLP enables to make good use of low-level and mid-level features that include useful intermediate-transition information [40]. CMLP utilizes MLP to encode input data into multiple dimensions [64, 128, 256, 512, 1024]. Then, we maxpool the output of the last four layers to construct a multiple-dimensional feature vector as follows:

$$\mathbf{v}_1 \in \mathbb{R}^{128}, \mathbf{v}_2 \in \mathbb{R}^{256}, \mathbf{v}_3 \in \mathbb{R}^{512}, \mathbf{v}_4 \in \mathbb{R}^{1024} \quad (30)$$

The combined feature vector is constructed as: $\bar{\mathbf{v}}_i = [\mathbf{v}_1^T, \mathbf{v}_2^T, \mathbf{v}_3^T, \mathbf{v}_4^T]^T \in \mathbb{R}^{1920}$. Three independent CMLPs map three resolutions into three individual $\bar{\mathbf{v}}_i$, for $i = 1, 2, 3$. Each $\bar{\mathbf{v}}_i$ represents the extracted potential information of each resolution. Then, the augmented feature matrix \mathbf{V} is generated by arranging its columns as $\mathbf{V} = [\bar{\mathbf{v}}_1, \bar{\mathbf{v}}_2, \bar{\mathbf{v}}_3] \in \mathbb{R}^{1920 \times 3}$, and further through 1D-convolution to obtain $\mathbf{M} \in \mathbb{R}^{N \times 3}$. Finally, the predicted next-moment shape $\hat{\mathbf{c}}_{k+1} \in \mathbb{R}^{3N}$ is obtained by vectorizing \mathbf{M} . The prediction loss of MRE is:

$$L_{mre} = \|\hat{\mathbf{c}}_{k+1} - \bar{\mathbf{c}}_{k+1}\| \quad (31)$$

where $\bar{\mathbf{c}}_{k+1}$ represents the ground-truth next-moment shape in the training data-set.

C. Discriminator Network

Generative Adversarial Network (GAN) is chosen as DN to enhance the prediction accuracy. For simplicity, we define $\Phi = \text{MRE}()$ and $\Psi = \text{DN}()$. We define $(\bar{\mathbf{c}}_k, \mathbf{r}_k, \mathbf{u}_k)$ as the \mathcal{X} input into Φ , while \mathcal{Y} represents the true shape $\bar{\mathbf{c}}_{k+1}$. Ψ is a classification network with similar structure as CMLP, constituted by serial MLP layers [128, 256, 512, 1024] to distinguish the predicted shape $\Phi(\mathcal{X})$ and the real shape \mathcal{Y} . We maxpool the last three layers of Ψ to obtain feature vector [256, 512, 1024]. Three feature vectors are concatenated into a latent vector $\mathbf{m} \in \mathbb{R}^{1792}$, and then passed through the fully-connected layers [512, 256, 64, 1] followed by sigmoid-classifier to obtain the evaluation. The adversarial loss is defined as follows:

$$L_{adv} = \sum_{1 \leq i \leq v} \log(1 - \Psi(\beta_i)) + \sum_{1 \leq j \leq v} \log(\Psi(\Phi(\alpha_j))) \quad (32)$$

where $\alpha_i \in \mathcal{X}, \beta_i \in \mathcal{Y}, i = 1, \dots, v$, and v is size of the dataset including \mathcal{X} and \mathcal{Y} . The total loss of SPN is:

$$L = \zeta_{mre} L_{mre} + \zeta_{adv} L_{adv} \quad (33)$$

where ζ_{mre} and ζ_{adv} are the weights of L_{mre} and L_{adv} , respectively, which satisfy the condition: $\zeta_{mre} + \zeta_{adv} = 1$.

V. RECEDING-TIME MODEL ESTIMATION

In this paper, the objects are assumed to be manipulated by the robot slowly, thus $\mathbf{J}_k(\mathbf{s}_k, \mathbf{r}_k)$ is expected to change smoothly. For ease of presentation, we omit the arguments of $\mathbf{J}_k(\mathbf{s}_k, \mathbf{r}_k)$ and denote it as \mathbf{J}_k from now on. To estimate the DJM, three indicators are considered, viz., accuracy, smoothness, and singularity. To this end, an optimization-based receding-time model (RTM) estimator is presented to estimate the changes of the Jacobian matrix, denoted by $\Delta \hat{\mathbf{J}}_k = \hat{\mathbf{J}}_k - \hat{\mathbf{J}}_{k-1}$, which enables to monitor the estimation procedure. $\Delta \hat{\mathbf{J}}_k$ can be obtained by considering the following three constraints:

- (Q_1) Constraint of receding-time error [43]. As \mathbf{J}_k depicts the relationship between $\Delta \mathbf{s}_{k+1}$ and \mathbf{u}_k in a local range, thus we consider the accumulated error in η past moments to ensure the estimation accuracy. η is the receding window size. The receding-time error is given by:

$$Q_1 = \sum_{j=1}^{\eta} \gamma^j \left\| \Delta \mathbf{s}_{k+1-j} - (\hat{\mathbf{J}}_{k-1} + \Delta \hat{\mathbf{J}}_k) \mathbf{u}_{k-j} \right\|^2 \quad (34)$$

The sensitivity to noise can be improved by adjusting η , which helps to address the measurement fluctuations. $0 < \gamma \leq 1$ is a constant forgetting factor giving less weight to the past observation data.

- (Q_2) Constraint of estimation smoothness [43]. As \mathbf{J}_k is assumed to be smooth, thus $\Delta \hat{\mathbf{J}}_k$ should be estimated smoothly to avoid sudden large fluctuations, which can be achieved by minimizing the Frobenius norm of $\Delta \hat{\mathbf{J}}_k$:

$$Q_2 = \|\Delta \hat{\mathbf{J}}_k\|_F^2 \quad (35)$$

- (Q_3) Constraint of shape manipulability [44]. It evaluates the feasibility of changing the object's shape under the

current object-robot configuration:

$$Q_3 = \left\| \frac{\lambda_{\max}((\hat{\mathbf{J}}_{k-1} + \Delta \hat{\mathbf{J}}_k)^\top (\hat{\mathbf{J}}_{k-1} + \Delta \hat{\mathbf{J}}_k))}{\lambda_{\min}((\hat{\mathbf{J}}_{k-1} + \Delta \hat{\mathbf{J}}_k)^\top (\hat{\mathbf{J}}_{k-1} + \Delta \hat{\mathbf{J}}_k))} \right\|^2 \quad (36)$$

where λ_{\max} and λ_{\min} are the maximum and minimum eigenvalue, respectively, and $Q_3 \geq 1$. When $Q_3 = 1$, the object can deform isotropically in any direction. A growing Q_3 indicates that the object is reaching singular (non-manipulable) configuration.

Finally, the total weighted optimization index is given by:

$$Q(\Delta \hat{\mathbf{J}}_k) = \mu_1 Q_1 + \mu_2 Q_2 + \mu_3 Q_3 \quad (37)$$

where $\mu_i > 0$ are the weights that specify the contribution of each constraint, and which satisfy $\mu_1 + \mu_2 + \mu_3 = 1$. The index (37) is then solved by using numerical optimization tools (e.g., *Matlab/fmincon* or *Python/CasADi*) to obtain $\Delta \hat{\mathbf{J}}_k$ and thus, iteratively update the deformation Jacobian matrix as follows: $\hat{\mathbf{J}}_k = \hat{\mathbf{J}}_{k-1} + \Delta \hat{\mathbf{J}}_k$.

VI. MODEL PREDICTIVE CONTROLLER

It is assumed that the matrix \mathbf{J}_k has been accurately estimated at the time instant k by the RTM, such that it satisfies $\Delta \mathbf{s}_{k+1} = \hat{\mathbf{J}}_k \mathbf{u}_k$. Based on this model, we propose an MPC-based controller to derive the velocity inputs \mathbf{u}_k for the robot, while taking saturation and workspace constraints into account. Two vectors are defined as follow:

$$\begin{aligned} \bar{\mathbf{s}}_k &= [\mathbf{s}_{k+1|k}^\top, \dots, \mathbf{s}_{k+h|k}^\top]^\top \in \mathbb{R}^{ph} \\ \bar{\mathbf{u}}_k &= [\mathbf{u}_{k|k}^\top, \dots, \mathbf{u}_{k+h-1|k}^\top]^\top \in \mathbb{R}^{3h} \end{aligned} \quad (38)$$

where $\bar{\mathbf{s}}_k$ and $\bar{\mathbf{u}}_k$ represent the predictions of \mathbf{s}_k and \mathbf{u}_k in the next h periods, respectively. The vectors $\mathbf{s}_{k+i|k}$ and $\mathbf{u}_{k+i|k}$ denote the i th predictions of \mathbf{s}_k and \mathbf{u}_k from the time instant k , where $\mathbf{s}_{k|k} = \mathbf{s}_k$, and $\mathbf{u}_{k|k} = \mathbf{u}_k$ must hold. The prediction $\bar{\mathbf{s}}_k$ can be calculated from the estimated Jacobian matrix by noting that $\hat{\mathbf{J}}_k \approx \hat{\mathbf{J}}_{k+h}$ is satisfied during period $[k, k+h]$ (which is reasonable, given the regularity of the object). This way, the predictions are computed as follows:

$$\mathbf{s}_{k+j|k} = \mathbf{s}_k + \sum_{i=0}^{j-1} \hat{\mathbf{J}}_k \mathbf{u}_{k+i|k}, \quad j = 1, \dots, h \quad (39)$$

All predictions are then grouped and arranged into a single vector form:

$$\begin{aligned} \bar{\mathbf{s}}_k &= \mathbf{A} \mathbf{s}_k + \Theta \bar{\mathbf{u}}_k, \\ \mathbf{A} &= \mathbf{I}_{h \times 1} \otimes \mathbf{E}_p \in \mathbb{R}^{ph \times p}, \quad \Theta = \mathbf{L}_h \otimes \hat{\mathbf{J}}_k \in \mathbb{R}^{ph \times 3h} \end{aligned} \quad (40)$$

In addition to $\bar{\mathbf{s}}_k$ and $\bar{\mathbf{u}}_k$, we define the constant sequence vector $\bar{\mathbf{s}}_k^*$ that represents the desired shape feature as:

$$\bar{\mathbf{s}}_k^* = [\mathbf{s}_{k+1|k}^{*\top}, \dots, \mathbf{s}_{k+h|k}^{*\top}]^\top \in \mathbb{R}^{ph} \quad (41)$$

The cost function $Q(\bar{\mathbf{u}}_k)$ for the optimization of the control input is formulated as:

$$Q(\bar{\mathbf{u}}_k) = (\bar{\mathbf{s}}_k - \bar{\mathbf{s}}_k^*)^\top \Upsilon_1 (\bar{\mathbf{s}}_k - \bar{\mathbf{s}}_k^*) + \bar{\mathbf{u}}_k^\top \Upsilon_2 \bar{\mathbf{u}}_k \quad (42)$$

where $\Upsilon_1 > 0$ and $\Upsilon_2 > 0$ are the weights for the error convergence rate and the smoothness of \mathbf{u}_k , respectively. Two constraints are considered:

- *Saturation limits.* In practice, robots have limits on their achievable joint speeds. These constraints are useful in soft object manipulation tasks to avoid damaging the object. Therefore, $\bar{\mathbf{u}}_k$ needs to be constrained:

$$\bar{\mathbf{u}}_{\min} \leq \bar{\mathbf{u}}_k \leq \bar{\mathbf{u}}_{\max} \quad (43)$$

where $\bar{\mathbf{u}}_{\min}$ and $\bar{\mathbf{u}}_{\max}$ are the constant lower and upper bounds, respectively.

- *Workspace limits.* Robots are also often required to operate in a confined workspace to avoid colliding with the environment. In soft object manipulation, this constraint is needed to avoid over-stretching or over-compressing the manipulated object. To this end, the following constant bounds are introduced:

$$\mathbf{r}_{k+i|k}^{\min} \leq \mathbf{r}_{k+i|k} \leq \mathbf{r}_{k+i|k}^{\max}, \quad 0 \leq i \leq h-1 \quad (44)$$

Similarly as in (39), the recursive structure of (44) can be obtained follows:

$$\mathbf{E}_{\min} \leq \mathbf{C}\bar{\mathbf{u}}_k \leq \mathbf{E}_{\max} \quad (45)$$

for $\mathbf{E}_{\min}, \mathbf{E}_{\max} \in \mathbb{R}^{3h}$ and $\mathbf{C} \in \mathbb{R}^{3h \times 3h}$ defined as:

$$\begin{aligned} \mathbf{E}_{\min} &= [(\mathbf{r}_{k|k}^{\min} - \mathbf{r}_{k-1})^\top, \dots, (\mathbf{r}_{k+h-1|k}^{\min} - \mathbf{r}_{k-1})^\top]^\top \\ \mathbf{E}_{\max} &= [(\mathbf{r}_{k|k}^{\max} - \mathbf{r}_{k-1})^\top, \dots, (\mathbf{r}_{k+h-1|k}^{\max} - \mathbf{r}_{k-1})^\top]^\top \\ \mathbf{C} &= \mathbf{L}_h \otimes \mathbf{E}_3 \in \mathbb{R}^{3h \times 3h} \end{aligned} \quad (46)$$

The quadratic optimization problem is formulated as follows:

$$\begin{aligned} \min_{\bar{\mathbf{u}}_k} Q(\bar{\mathbf{u}}_k) &= \frac{1}{2} \bar{\mathbf{u}}_k^\top \mathbf{H} \bar{\mathbf{u}}_k + \mathbf{q}^\top \bar{\mathbf{u}}_k \\ \text{s.t.} \quad &\bar{\mathbf{u}}_{\min} \leq \bar{\mathbf{u}}_k \leq \bar{\mathbf{u}}_{\max} \\ &\mathbf{E}_{\min} \leq \mathbf{C}\bar{\mathbf{u}}_k \leq \mathbf{E}_{\max} \end{aligned} \quad (47)$$

where $\mathbf{H} = 2(\Theta^\top \Upsilon_1 \Theta + \Upsilon_2) \in \mathbb{R}^{3h \times 3h}$, $\mathbf{q}^\top = 2\Omega^\top \Upsilon_1 \Theta \in \mathbb{R}^{1 \times 3h}$ and $\Omega = \mathbf{A}\mathbf{s}_k - \bar{\mathbf{s}}_k^* \in \mathbb{R}^{ph}$ are constant matrices. Then, $\bar{\mathbf{u}}_k$ can be obtained by using a standard quadratic solver on (47). Finally, \mathbf{u}_k is calculated by the receding horizon scheme:

$$\mathbf{u}_k = [\mathbf{E}_3, \mathbf{0}, \dots, \mathbf{0}] \bar{\mathbf{u}}_k \quad (48)$$

Fig. 5 presents the conceptual block diagram of the proposed framework.

Remark 2. The proposed MPC-based technique (47) computes the robot's shaping actions based on a performance objective and subject to system's constraints; This approach does not require the identification of the full analytical model of the deformable object. Its quadratic optimization form enables to integrate additional metrics into the problem, e.g., rising time and overshoot.

VII. RESULTS

A. Experimental Setup

Vision-based manipulation experiments are conducted to validate our proposed framework. The experimental platform used in our study includes a fixed D455 depth sensor, a UR5 robot manipulator, and various deformable objects shown in Fig. 6. The depth sensor receives the video stream, from which it computes the 3D shapes by using the OpenCV and RealSense libraries. In our experiments, only 3 DOF of the robot manipulator are considered, therefore, the control

input $\mathbf{u} = [u_x, u_y, u_z]^\top \in \mathbb{R}^3$ represents the linear velocity of the end-effector; A saturation limit of $|u_i| \leq 0.01$ m/s, is applied for $i = x, y, z$. The motion control algorithm is implemented on ROS/Python, which runs with a servo-control loop of 10 Hz. A video of the conducted experiments can be downloaded from https://github.com/JiamingQi-Tom/experiment_video/raw/master/paper4/video.mp4

The proposed shape extraction algorithm is depicted in Fig. 7. The RGB image from the camera is transformed into HSV and combined with mask processing to obtain a binary image. *OpenCV/thinning* is utilized with FPS to extract a fixed number of object points. The point on the centerline closest to the gripper's green marker is chosen to sort the centerline along the cable. Then we obtain the 3D shapes by using the RealSense sensor and checking its 2D pixels. We adopt the method in [11] to compute the object's contour. A surface is obtained in the similar way to the centerline, i.e., by sorting points from top to bottom, and from left to right. As 2D pixels and 3D points have a one-to-one correspondence in a depth camera, thus, our extraction method improves the robustness to measurement noise and is simpler than traditional point cloud processing algorithms.

TABLE I
FITTING CONFIGURATIONS. "N/A" STANDS FOR "NOT APPLICABLE".

	Centerline	Contour	Surface / Rigid
Method	LSM	MLS	MLS
Support Radius	N/A	$d = 0.2$	$d = 0.2$
PCA	N/A	$m = 1$	$m = 1$
Fitting order	$n = 5$	$n = 4$	$n_x = n_y = 2$
Basis-Function	Bernstein	Trigonometric	Polynomial
Numbers	$N = 64$	$N = 64$	$N = 32$

B. Online Fitting of the Parametric Shape Representation

In this section, ten thousand samples of centerlines, contours and surfaces with $N = 64, 64, 32$, respectively, are collected by commanding the robot to manipulate the objects, whose configuration is then captured by a depth sensor. Such shaping actions are shown in Fig. 8, and visualized in the accompanying multimedia attachment. This data is used to evaluate the performance (viz. its accuracy and computation time) of our representation framework. For that, we calculate the average error between the feedback shape $\bar{\mathbf{c}}$ and the reconstructed shape $\hat{\mathbf{c}}$ as follows $\text{mean}(\sum \|\bar{\mathbf{c}}_i - \hat{\mathbf{c}}_i\|)$. The computation time is defined as the average of the overall processing time of all sample data among each method.

Fig. 9 shows that the larger the scalars n, n_x, n_y are, the better the fitting accuracy of LSM and MLS is. MLS fits better than LSM under the same condition, as MLS calculates the independent weight while LSM assumes that each node has the same weight. MLS works better in fitting contour because the parametric curve may not be continuous in the end corner, thus, the equal weight assumption of LSM is not suitable here. As the number of data points for surface is $N = 32$, it does not satisfy the condition $N \gg (n_x + 1)(n_y + 1)$ for higher order fitting models (e.g., $n_x = n_y \geq 5$), thus, we only use $n_x = n_y \leq 4$. The results show that MLS performs better than LSM in the surface representation; Interestingly, MLS

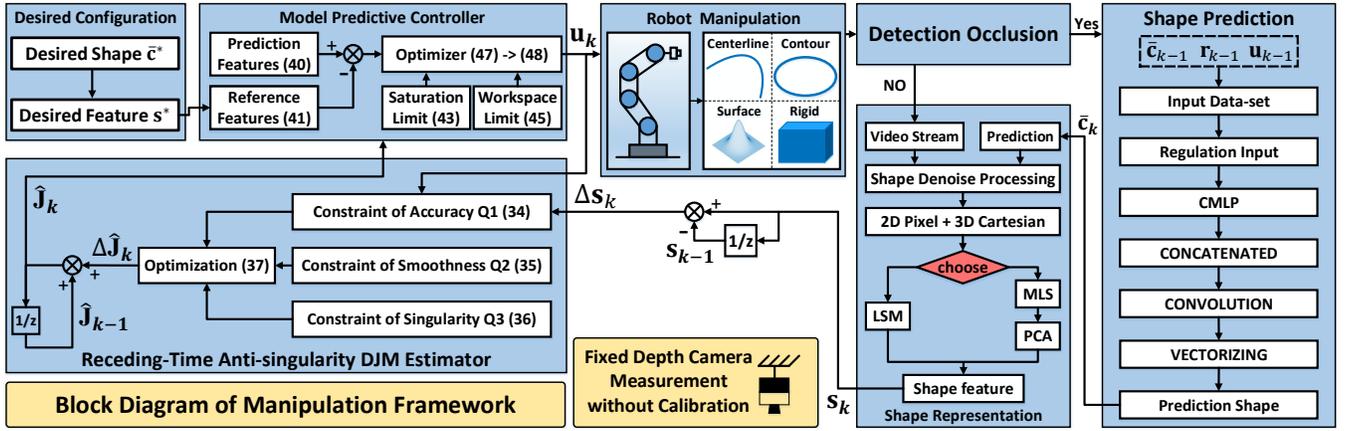


Fig. 5. The block diagram of the proposed shape servoing framework, including representation, prediction, approximation, and manipulation within constraints.

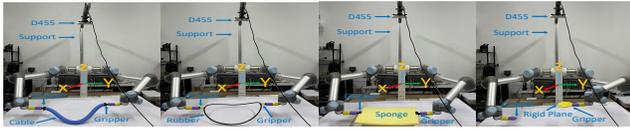


Fig. 6. The experimental setup, including the objects (elastic and rigid), single-arm robot (UR5), and D455.

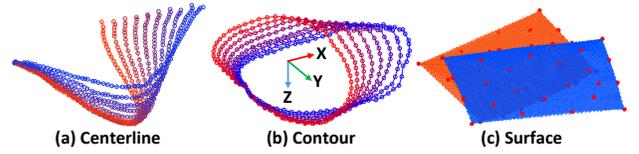


Fig. 8. Shapes of various objects manipulated by UR5.

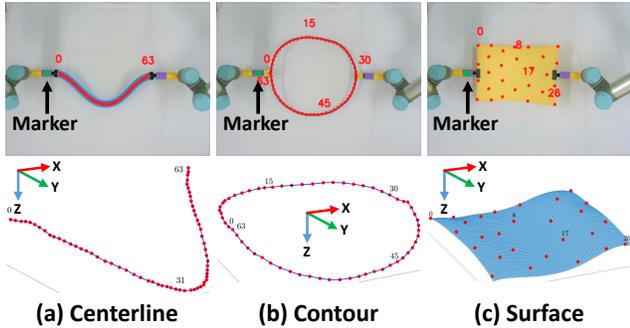


Fig. 7. Shape extraction for centerline, contour and surface. The first row shows the 2-D image pixels, and the second row shows the 3-D shapes. All shapes are fixed-sampled, equidistant, and ordered sorting.

obtains satisfactory performance even with $n_x = n_y = 1$. Fig. 9 also shows that larger n, n_x, n_y will also increase the computation time. MLS has a more noticeable increase, as it calculates the weights of all nodes while LSM calculates them once. The trigonometric approach is the fastest, polynomial and Bernstein follow, while Cox-deBoor is the slowest with the most iterative operations. The above analysis verifies the effectiveness of the proposed extraction framework, which can represent objects with a low-dimensional feature. Details of the curve fitting configuration is given in Table. I.

C. Occlusion-Robust Prediction of Object Shapes

The data collected in Section VII-B is used to evaluate our proposed SPN. 80% of the data is used as the training set, and the remaining 20% of the data is used to test the trained network. We set $\delta = 2$ for the centerline, contour, and surface parametrization. SPN is built using PyTorch and trained by an ADAM optimizer with a batch size of 500, and the initial learning rate set to 0.0001. RELU activation and batch normalization are adopted to improve the network's

performance. In this validation test, the robot deforms the objects with small babbling motions while a cardboard sheet covers parts of objects. Fig. 10 shows that SPN can predict and provide relatively complete shapes for the three types of manipulated objects. The accompanying video demonstrates the performance of this method.

D. Estimation of the Sensorimotor Model

This section aims to evaluate RTM (37) that approximates the deformation Jacobian matrix; The performance of the RTM estimator is compared with the interaction matrix estimator (IEM) in [15], and the unscented Kalman filter (UKF) in [31]. To this end, we introduce two metrics, i.e., T_1 and T_2 , to quantitatively compare performances of these algorithms:

$$T_1 = \|\hat{\mathbf{s}}_{k+1} - \mathbf{s}_{k+1}\|, \quad T_2 = \|\Delta \mathbf{s}_{k+1} - \hat{\mathbf{J}}_k \mathbf{u}_k\| \quad (49)$$

where $\hat{\mathbf{s}}_{k+1} = \hat{\mathbf{s}}_k + \hat{\mathbf{J}}_k \mathbf{u}_k$ is the approximated shape feature that is computed based on the control actions. Fig. 11 shows that RTM with $\eta = 20$ provides the best performance for T_1 among the methods; This means that constraint Q_1 (34) enables RTM to learn from past data by adjusting η . Small values of T_2 reflect that RTM accurately predicts the differential changes induced by the DJM; This means that constraint Q_2 (35) helps to compute a smooth matrix $\hat{\mathbf{J}}_k$. As the proposed method incorporates the constraint Q_3 (36), thus, RTM can prevent singularities in the estimation of the Jacobian matrix, while IEM and UKF are prone to reach ill-conditioned estimations. We use RTM with $\eta = 10$ in the following sections.

E. Automatic Shape Servoing Control

This section conducts four automatic manipulation experiments, labeled as Exp1, Exp2, Exp3, and Exp4, respectively. The target contour $\bar{\mathbf{c}}^*$ is obtained from previous demonstrations of the manipulation task, which ensures its reachability.

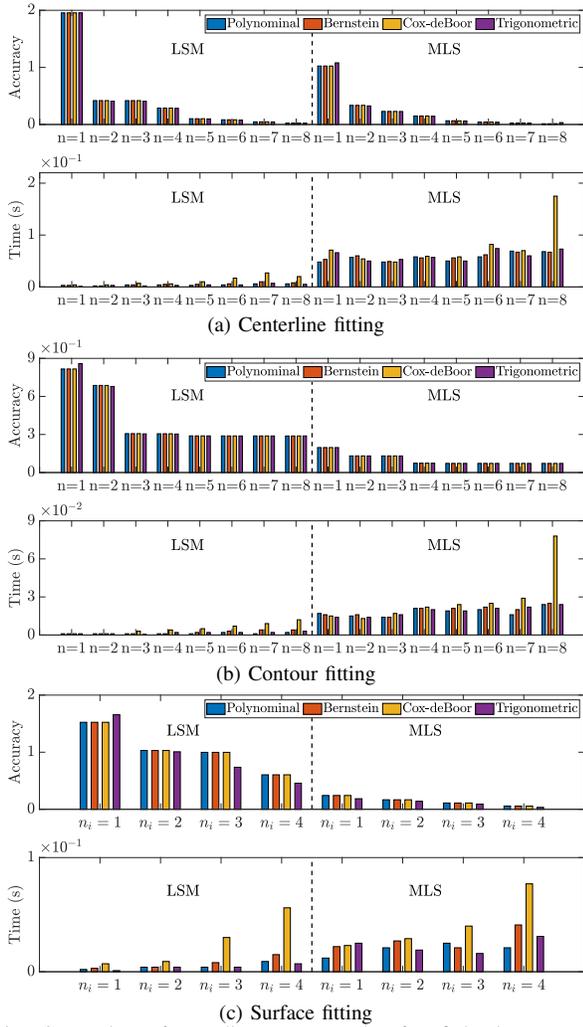


Fig. 9. Comparison of centerline, contour and surface fitting in accuracy and time-consuming among (7)(8)(9)(10) between LSM and MLS with $d = 0.9$, $d = 0.1$, and $d = 0.2$, respectively. $n_i, i = x, y$ are the surface fitting order.

A cardboard sheet is manually placed over the object to produce (partial) occlusions and test the robustness of our algorithm. The estimation methods in [15] and [31] are compared with the proposed receding-time model with MPC ($h = 5$ and $h = 15$). We label these methods as CM_1, \dots, CM_4 , and each method has been optimized to achieve a balance of stability, convergence, and responsiveness.

In addition to the feedback shape error $\|s^* - s_k\|$, we also compare T_{max} (i.e., the number of steps from start to finish), t_d (the steps from $\|s^* - s_0\|$ to 10% of this value), t_s (the steps from 10% $\|s^* - s_0\|$ to the threshold value), and d_{eff} (the total moving distance of the end-effector), [11], [23]. Fig. 12 shows the shaping motions of the manipulated objects toward the desired configuration (black curve), with the cardboard blocking the view at various instances. The results demonstrate that SPN can predict the object's shape during occlusions and feed it back to the controller to enforce the shape servo-loop; This results in the objects being gradually manipulated towards the desired configuration. The error norm $\|s^* - s_k\|$ plots in Fig. 13 shows that CM_3 provides the best control performance for the error minimization, with CM_4 as the second-best, and CM_1 and CM_2 showing similar performance.

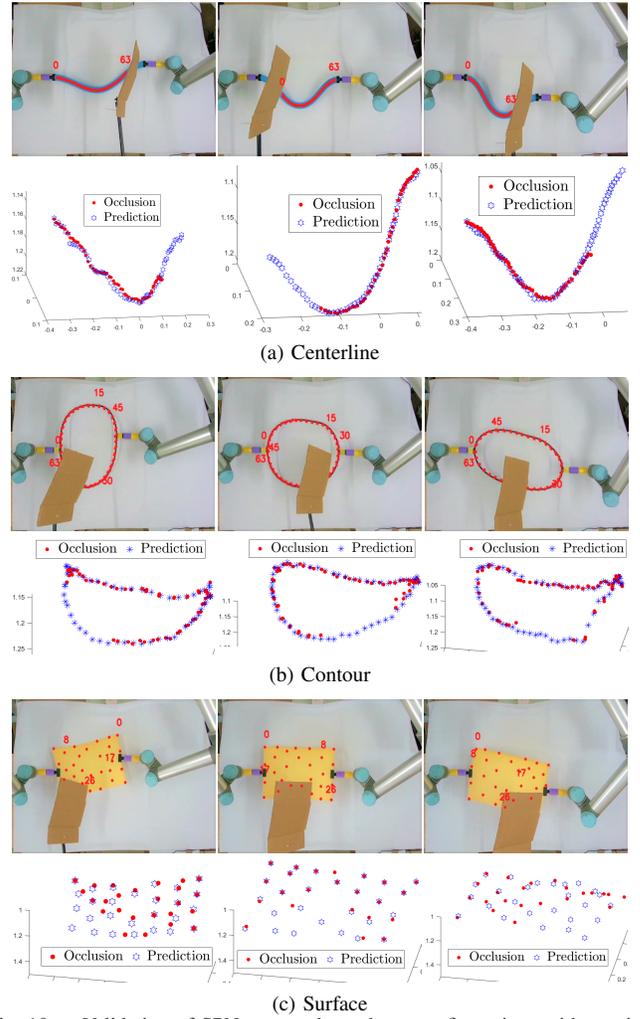


Fig. 10. Validation of SPN among three shape configurations with moving the obstacle manually. The RGB image describes the occluded case. Red dot in the second row gives the feedback shape within the occlusion at the current instant, and blue ones are the predicted shape according to the past data.

A higher h is helpful for feature prediction, yet, since we assume that \hat{J}_k is constant in the window period h , it may lead to inaccurate predictions and even wrong manipulation of objects. Therefore, h should be chosen according to the performance requirements of the system.

From the Q_3 plots in Fig. 13, we can see RTM with $\eta = 10$ provides the smallest value, which validates that RTM can enhance the manipulation feasibility and avoid shapes falling into the singular configurations. The black dashed lines represent the workspace constraints of r_x, r_y, r_z , on which we can see the end-effector's trajectories in the Cartesian coordinate system; The results show that CM_3 and CM_4 remain within the workspace, while CM_1 and CM_2 may violate this constraint. Due to the adopted input saturation in our platform, all four methods can satisfy the control saturation constraint. Exp4 shows that our method has good universality, not only for the shape control, but also for traditional rigid object positioning.

A performance comparison of these experiments is given in Fig. 14. The results illustrate that our proposed shape servoing framework achieves the best performance relative to

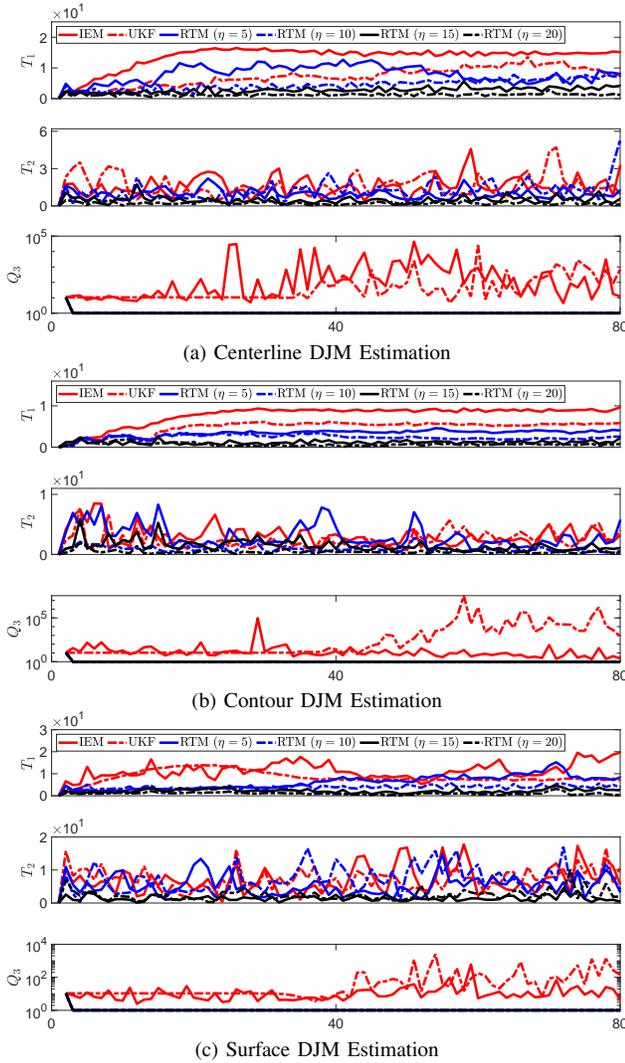


Fig. 11. Curves of T_1 , T_2 , and Q_3 , reviewing the accuracy, smoothness, and singularity of DJM, respectively. T_1 represents the feature estimation error, and T_2 depicts the estimation differential error. Besides, $\mu_1 = 0.8$, $\mu_2 = 0.1$, $\mu_3 = 0.1$.

the manipulation speed T_{max} , response speed t_d , convergence speed t_s , and motion distance d_{eff} .

VIII. CONCLUSIONS

In this paper, we present an occlusion-robust shape servoing framework to control shapes of elastic objects into target configurations, while considering workspace and saturation constraints. A low-dimensional feature extractor is proposed to represent 3D shapes based on LSM and MLS. A deep neural network is introduced to predict the object's configuration subject to occlusions, and feed it to the shape servo-controller. A receding-time model estimator is designed to approximate the deformation Jacobian matrix with various constraints such as accuracy, smoothness, singularity. The conducted experiments validate the proposed methodology with multiple unstructured shape servoing tasks in visually occluded situations and with unknown deformation models.

However, there are some limitations in our framework. For example, as the support field radius d is constant (i.e., it does not adjust with dynamic shapes), the computed representation

lacks flexibility. Also, the SPN needs to obtain substantial offline training data to properly work, which might pose complications in practice. Note that our method may not accurately shape objects with negligible elastic properties (e.g., fabrics, food materials, etc). Future work include the incorporation of shape reachability detection into the framework in order to determine the feasibility of a given shaping task beforehand.

REFERENCES

- [1] J. Zhu, A. Cherubini, C. Dune, D. Navarro-Alarcon, F. Alambeigi, D. Berenson, F. Ficuciello, K. Harada, X. Li, J. Pan, *et al.*, "Challenges and outlook in robotic manipulation of deformable objects," *arXiv preprint arXiv:2105.01767*, 2021.
- [2] X. Li, X. Su, and Y.-H. Liu, "Vision-based robotic manipulation of flexible pcbs," *IEEE/ASME Transactions on Mechatronics*, vol. 23, no. 6, pp. 2739–2749, 2018.
- [3] L. Cao, X. Li, P. T. Phan, A. M. H. Tiong, H. L. Kaan, J. Liu, W. Lai, Y. Huang, H. M. Le, M. Miyasaka, *et al.*, "Sewing up the wounds: A robotic suturing system for flexible endoscopy," *IEEE Robotics & Automation Magazine*, vol. 27, no. 3, pp. 45–54, 2020.
- [4] Z. Hu, P. Sun, and J. Pan, "Three-dimensional deformable object manipulation using fast online gaussian process regression," *IEEE Robotics and Automation Letters*, vol. 3, no. 2, pp. 979–986, 2018.
- [5] J. Sanchez, J.-A. Corrales, B.-C. Bouzgarrou, and Y. Mezouar, "Robotic manipulation and sensing of deformable objects in domestic and industrial applications: a survey," *Int. J. of Rob. Res.*, 2018.
- [6] H. Yin, A. Varava, and D. Kragic, "Modeling, learning, perception, and control methods for deformable object manipulation," *Science Robotics*, vol. 6, no. 54, 2021.
- [7] D. Navarro-Alarcon, Y.-H. Liu, J. G. Romero, and P. Li, "Model-free visually servoed deformation control of elastic objects by robot manipulators," *IEEE Transactions on Robotics*, vol. 29, no. 6, pp. 1457–1468, 2013.
- [8] D. Navarro-Alarcon, Y.-h. Liu, J. G. Romero, and P. Li, "On the visual deformation servoing of compliant objects: Uncalibrated control methods and experiments," *The International Journal of Robotics Research*, vol. 33, no. 11, pp. 1462–1480, 2014.
- [9] H. Wang, B. Yang, J. Wang, X. Liang, W. Chen, and Y.-H. Liu, "Adaptive visual servoing of contour features," *IEEE/ASME Transactions on Mechatronics*, vol. 23, no. 2, pp. 811–822, 2018.
- [10] D. Navarro-Alarcon and Y.-H. Liu, "Fourier-based shape servoing: a new feedback method to actively deform soft objects into desired 2-d image contours," *IEEE Transactions on Robotics*, vol. 34, no. 1, pp. 272–279, 2017.
- [11] J. Qi, G. Ma, J. Zhu, P. Zhou, Y. Lyu, H. Zhang, and D. Navarro-Alarcon, "Contour moments based manipulation of composite rigid-deformable objects with finite time model estimation and shape/position control," *IEEE/ASME Transactions on Mechatronics*, 2021.
- [12] Z. Hu, T. Han, P. Sun, J. Pan, and D. Manocha, "3-d deformable object manipulation using deep neural networks," *IEEE Robotics and Automation Letters*, vol. 4, no. 4, pp. 4255–4261, 2019.
- [13] J. Qi, G. Ma, P. Zhou, H. Zhang, Y. Lyu, and D. Navarro-Alarcon, "Towards latent space based manipulation of elastic rods using autoencoder models and robust centerline extractions," *Advanced Robotics*, vol. 36, no. 3, pp. 101–115, 2022.
- [14] P. Zhou, J. Zhu, S. Huo, and D. Navarro-Alarcon, "LaSeSOM: A latent and semantic representation framework for soft object manipulation," *IEEE Robotics and Automation Letters*, vol. 6, no. 3, pp. 5381–5388, 2021.
- [15] J. Zhu, D. Navarro-Alarcon, R. Passama, and A. Cherubini, "Vision-based manipulation of deformable and rigid objects using subspace projections of 2d contours," *Robotics and Autonomous Systems*, vol. 142, p. 103798, 2021.
- [16] N. Cazy, P.-B. Wieber, P. R. Giordano, and F. Chaumette, "Visual servoing when visual information is missing: Experimental comparison of visual feature prediction schemes," in *2015 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2015, pp. 6031–6036.
- [17] T. Tang, C. Wang, and M. Tomizuka, "A framework for manipulating deformable linear objects by coherent point drift," *IEEE Robotics and Automation Letters*, vol. 3, no. 4, pp. 3426–3433, 2018.
- [18] T. Tang and M. Tomizuka, "Track deformable objects from point clouds with structure preserved registration," *The International Journal of Robotics Research*, p. 0278364919841431, 2018.

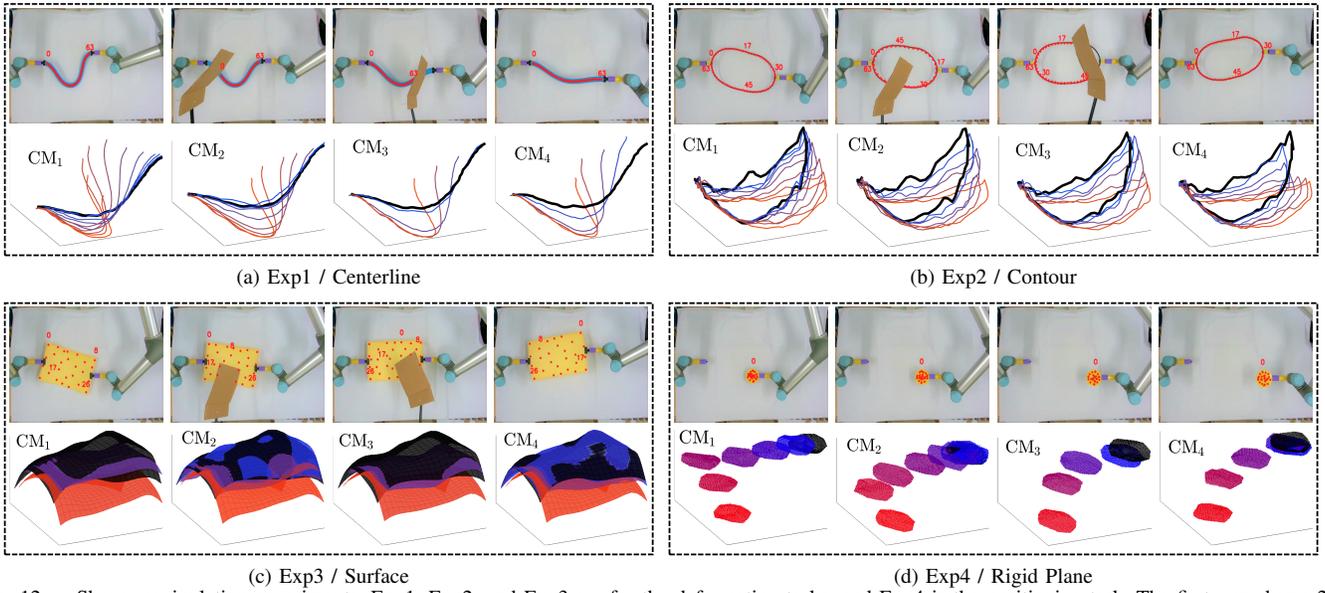


Fig. 12. Shape manipulation experiments, Exp1, Exp2, and Exp3 are for the deformation tasks, and Exp4 is the positioning task. The first row shows 2D image manipulation process (from left to right are initial, intermediate, intermediate, and desired shape), and the second row represents 3D shape manipulation process. The obstacle moves randomly during the process. The target shape is shown with black, the red represents the start shapes, and the gradient color are intermediate shapes.

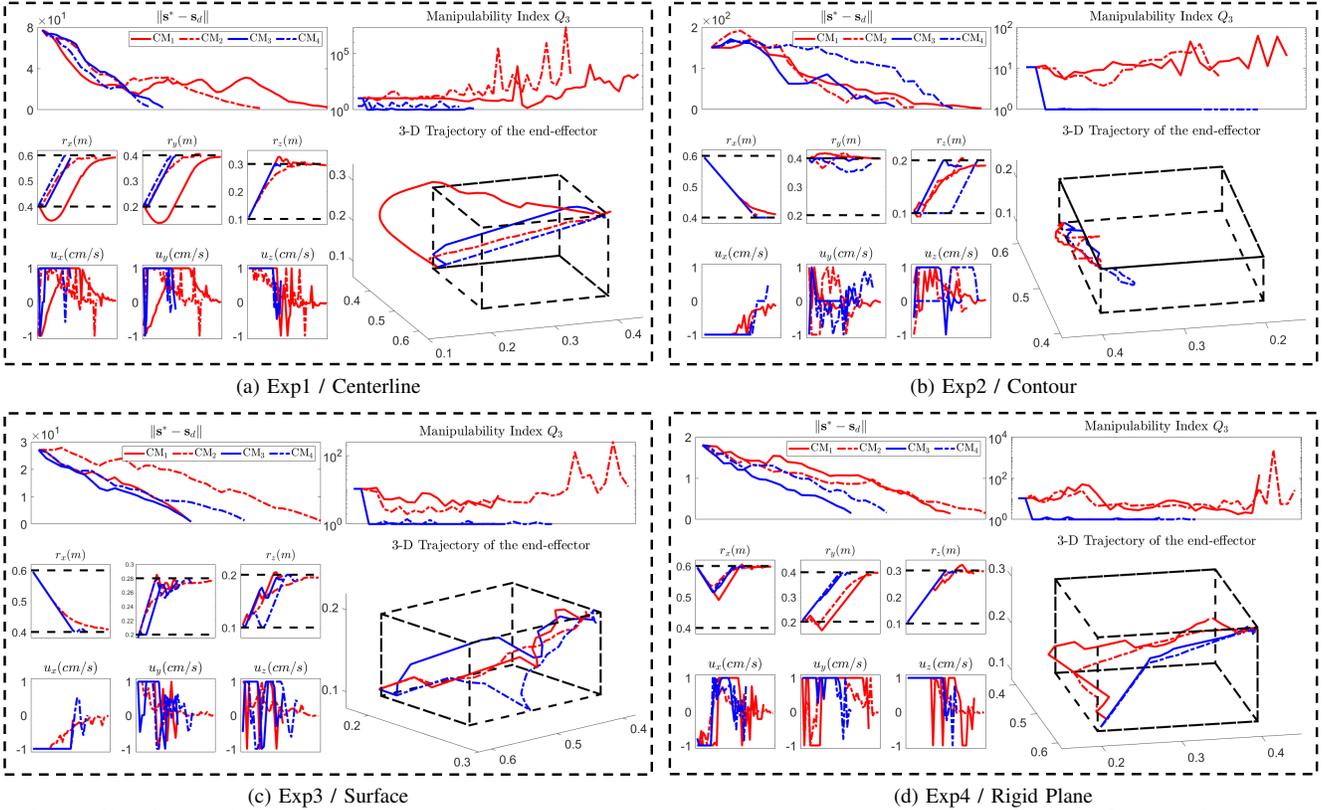


Fig. 13. Profiles of the manipulation error $\|s^* - s_k\|$, the robot pose \mathbf{r}_k , the velocity command \mathbf{u}_k , and the manipulability index Q_3 among four experiments (Exp1, ..., Exp4). Each experiment adopts four methods, i.e., CM_1, \dots, CM_4 .

- [19] D. Navarro-Alarcon, J. Qi, J. Zhu, and A. Cherubini, "A Lyapunov-stable adaptive method to approximate sensorimotor models for sensor-based control," *Frontiers in Neurorobotics*, vol. 14, 2020.
- [20] F. Chaumette and S. Hutchinson, "Visual servo control, part I: Basic approaches," *IEEE Robotics and Automation Magazine*, vol. 13, no. 4, pp. 82–90, 2006.
- [21] K. Hosoda and M. Asada, "Versatile visual servoing without knowledge of true jacobian," in *Proceedings of IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS'94)*, vol. 1. IEEE, 1994, pp. 186–193.
- [22] F. Alamebeigi, Z. Wang, R. Hegeman, Y.-H. Liu, and M. Armand, "Autonomous data-driven manipulation of unknown anisotropic deformable tissues using unmodelled continuum manipulators," *IEEE Robotics and Automation Letters*, vol. 4, no. 2, pp. 254–261, 2018.
- [23] R. Lagneau, A. Krupa, and M. Marchal, "Active deformation through visual servoing of soft objects," in *2020 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2020, pp. 8978–8984.
- [24] Lagneau, K. Romain, M. Alexandre, and Maud, "Automatic shape control of deformable wires based on model-free visual servoing," *IEEE Robotics and Automation Letters*, vol. 5, no. 4, pp. 5252–5259, 2020.

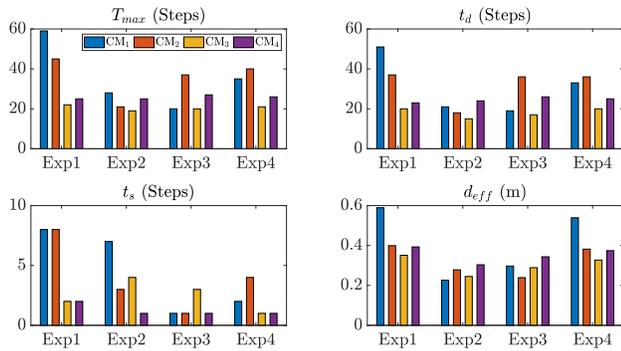


Fig. 14. Performance comparison in the manipulation tasks (Exp1 to Exp4).

- [25] L. Han, H. Wang, Z. Liu, W. Chen, and X. Zhang, "Visual tracking control of deformable objects with a fat-based controller," *IEEE Transactions on Industrial Electronics*, 2021.
- [26] A. Hajiloo, M. Keshmiri, W.-F. Xie, and T.-T. Wang, "Robust online model predictive control for a constrained image-based visual servoing," *IEEE Transactions on Industrial Electronics*, vol. 63, no. 4, pp. 2242–2250, 2015.
- [27] H. Wakamatsu and S. Hirai, "Static modeling of linear object deformation based on differential geometry," *The International Journal of Robotics Research*, vol. 23, no. 3, pp. 293–311, 2004.
- [28] M. Yu, H. Zhong, and X. Li, "Shape control of deformable linear objects with offline and online learning of local linear deformation models," *arXiv preprint arXiv:2109.11091*, 2021.
- [29] H. Abdi *et al.*, "The method of least squares," *Encyclopedia of Measurement and Statistics*. CA, USA: Thousand Oaks, 2007.
- [30] P. Lancaster and K. Salkauskas, "Surfaces generated by moving least squares methods," *Mathematics of computation*, vol. 37, no. 155, pp. 141–158, 1981.
- [31] J. Qi, W. Ma, D. Navarro-Alarcon, H. Gao, and G. Ma, "Adaptive shape servoing of elastic rods using parameterized regression features and auto-tuning motion controls," *arXiv preprint arXiv:2008.06896*, 2020.
- [32] K. Smith, "On the standard deviations of adjusted and interpolated values of an observed polynomial function and its constants and the guidance they give towards a proper choice of the distribution of observations," *Biometrika*, vol. 12, no. 1/2, pp. 1–85, 1918.
- [33] G. G. Lorentz, *Bernstein polynomials*. American Mathematical Soc., 2013.
- [34] I. J. Schoenberg, *Cardinal spline interpolation*. SIAM, 1973.
- [35] M. J. D. Powell *et al.*, *Approximation theory and methods*. Cambridge university press, 1981.
- [36] W. Y. Tey, N. A. Che Sidik, Y. Asako, M. W. Muhiedeen, and O. Afshar, "Moving least squares method and its improvement: A concise review," *Journal of Applied and Computational Mechanics*, vol. 7, no. 2, pp. 883–889, 2021.
- [37] H. Zhang, C. Guo, X. Su, and C. Zhu, "Measurement data fitting based on moving least squares method," *Mathematical Problems in Engineering*, vol. 2015, 2015.
- [38] J. C. Mason and D. C. Handscomb, *Chebyshev polynomials*. CRC press, 2002.
- [39] R. T. Farouki, "Legendre–bernstein basis transformations," *Journal of Computational and Applied Mathematics*, vol. 119, no. 1-2, pp. 145–160, 2000.
- [40] Z. Huang, Y. Yu, J. Xu, F. Ni, and X. Le, "Pf-net: Point fractal network for 3d point cloud completion," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 7662–7670.
- [41] X. Yan, C. Zheng, Z. Li, S. Wang, and S. Cui, "Pointasnl: Robust point clouds processing using nonlocal neural networks with adaptive sampling," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 5589–5598.
- [42] C. R. Qi, H. Su, K. Mo, and L. J. Guibas, "Pointnet: Deep learning on point sets for 3d classification and segmentation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 652–660.
- [43] H. Mo, B. Ouyang, L. Xing, D. Dong, Y. Liu, and D. Sun, "Automated 3-d deformation of a soft object using a continuum robot," *IEEE Transactions on Automation Science and Engineering*, 2020.
- [44] K. M. Lynch and F. C. Park, *Modern robotics*. Cambridge University Press, 2017.



Jiaming Qi received the M.Sc. in integrated circuit engineering from Harbin Institute of Technology, Harbin, China, in 2018. In 2019, he was a visiting PhD student at The Hong Kong Polytechnic University. He is currently pursuing the Ph.D. degree with control science and engineering, Harbin Institute of Technology, Harbin, China. His current research interests include data-driven control for soft object manipulation, visual servoing, robotics and control theory.



research interests include networked system cooperation, adaptive systems, and robotic control.

Dongyu Li (Member, IEEE) received the B.S. and Ph.D. degree in control science and engineering, Harbin Institute of Technology, China, in 2016 and 2020. He was a joint Ph.D. student with the Department of Electrical and Computer Engineering, National University of Singapore from 2017 to 2019, and a research fellow with the Department of Biomedical Engineering, National University of Singapore, from 2019 to 2021. He is currently an Associate Professor with the School of Cyber Science and Technology, Beihang University, China. His



Yufeng Gao (Student Member, IEEE) received his bachelor of engineering and bachelor of economics degrees from Wuhan University of Technology and Wuhan University, China, both in 2016. And he received the Ph.D. degree in control science and engineering, Harbin Institute of Technology, China, in 2021. He was a joint Ph.D. student with the Chair of Automatic Control Engineering, Technical University of Munich, Germany, from 2018 to 2019. His current research interests include spacecraft system modeling, control and optimization.



Peng Zhou (Student Member, IEEE) was born in China. He received the M.Sc. degree in software engineering from the School of Software Engineering, Tongji University, Shanghai, China, in 2017 and is currently pursuing his the Ph.D. degree in mechanical engineering at The Hong Kong Polytechnic University, Kowloon, Hong Kong. His research interests include deformable object manipulation, motion planning and robot learning.



Robotics and Machine Intelligence Laboratory. His current research interests include perceptual robotics and control theory.

David Navarro-Alarcon (Senior Member, IEEE) received the Ph.D. degree in mechanical and automation engineering from The Chinese University of Hong Kong (CUHK), Shatin, Hong Kong, in 2014.

From 2014 to 2017, he was a Postdoctoral Fellow and then a Research Assistant Professor with the CUHK T Stone Robotics Institute. Since 2017, he has been with The Hong Kong Polytechnic University, Hong Kong, where he is currently an Assistant Professor with the Department of Mechanical Engineering, and the Principal Investigator of the