

# ViT-VS: On the Applicability of Pretrained Vision Transformer Features for Generalizable Visual Servoing

Alessandro Scherl<sup>1,2</sup>, Stefan Thalhammer<sup>2</sup>, Bernhard Neuberger<sup>2</sup>, Wilfried Wöber<sup>2</sup>, José García-Rodríguez<sup>1</sup>

**Abstract**— Visual servoing enables robots to precisely position their end-effector relative to a target object. While classical methods rely on hand-crafted features and thus are universally applicable without task-specific training, they often struggle with occlusions and environmental variations, whereas learning-based approaches improve robustness but typically require extensive training. We present a visual servoing approach that leverages pretrained vision transformers for semantic feature extraction, combining the advantages of both paradigms while also being able to generalize beyond the provided sample. Our approach achieves full convergence in unperturbed scenarios and surpasses classical image-based visual servoing by up to 31.2% relative improvement in perturbed scenarios. Even the convergence rates of learning-based methods are matched despite requiring no task- or object-specific training. Real-world evaluations confirm robust performance in end-effector positioning, industrial box manipulation, and grasping of unseen objects using only a reference from the same category. Our code and simulation environment are available at: <https://alessandroscherl.github.io/ViT-VS/>

## I. INTRODUCTION

Visual servoing (VS) as a visual control strategy allows positioning the robot relative to a target with a single reference [1], [2].

This enables executing downstream tasks, such as object tracking and grasping [3], [4], [5], [6]. Generally, VS can be categorized into two approaches: Position-Based Visual Servoing (PBVS) operating on pose differences, and Image-Based Visual Servoing (IBVS), directly utilizing image features. These features range from geometric primitives and image moments [2] to feature descriptors [7] or direct utilization of photometric image data [8].

While effective in controlled settings, these classical methods show limited robustness to image perturbations and typically require exact target object instances. Recent learning-based approaches attempt to address these limitations through different strategies such as pose regression [9], velocity regression [10], learned feature extraction [11], metric learning [12], and unsupervised feature learning [13], [14], [5], [15]. However, these methods introduce new challenges - they require task-specific training, extensive data

\*This work has been submitted to the IEEE for possible publication. Copyright may be transferred without notice, after which this version may no longer be accessible.

<sup>1</sup>The authors are with the Department of Computer Technology, University of Alicante, Spain. (email: as358@alu.ua.es, jgarcia@dtic.ua.es)

<sup>2</sup>The authors are with the Industrial Engineering Department, UAS Technikum Vienna, Austria. (email: {alessandro.scherl, stefan.thalhammer, bernhard.neuberger, wilfried.woeber}@technikum-wien.at)

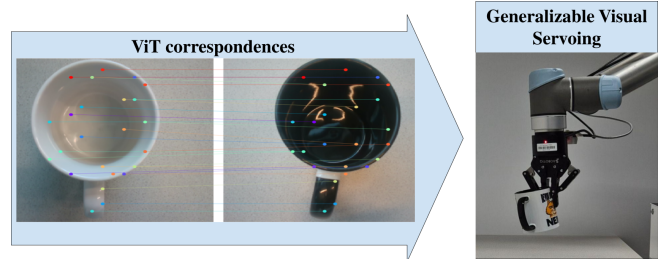


Fig. 1. ViT-VS category-level object grasping Left: ViT Correspondence Matching process with white mug as desired image. Right: Successful grasp of the target object.

generation, or predefined object models, making them difficult to deploy in real-world scenarios where unseen objects and environmental changes are common.

We hypothesize that pretrained Vision Transformers (ViTs) [16] effectively combine the advantages of both classical and learning-based VS approaches through semantically robust features. To validate this, we introduce ViT-VS, a generalizable visual servoing framework combining IBVS with DINOv2 [17] features. Our approach requires no task-specific training or fine-tuning, and is capable of performing visual servoing robust to image perturbations with high convergence rates. Fig. 1 showcases the ViT feature matching between different object instances (left) and the corresponding successful grasps (right). A key challenge in utilizing ViTs for VS stems from their rotation invariance property, developed during training for image classification [18]. This invariance can cause convergence to incorrect orientations, specifically with in-plane rotations of  $\pm 90^\circ$  and  $180^\circ$ . We address this by aligning the camera based on accumulated feature similarities across different simulated rotations of the current image before initiating the VS control loop. Additionally, the high computational complexity of ViT image processing leads to suboptimal path lengths during convergence. To alleviate this behavior we propose a velocity stabilization using an exponential moving average filter. We evaluate ViT-VS in both simulation and real-world environments, demonstrating:

- A novel VS approach that is applicable without training or fine-tuning, like classical IBVS approaches, yet achieves the convergence rates and robustness to image perturbations of learning-based approaches. Relative convergence rate improvement is 31.2% as compared to the best classical IBVS, achieving a 100% convergence rate for unperturbed scenarios.
- Strategies for addressing ViT limitations including compensation for ViTs' rotation invariance and trajectory

regularization for reducing convergence length ratios.

- Industrial box manipulation with 100% success rate ( $n = 20$ ) on a mobile robot with a starting point positioning error of  $\pm 10cm$ .
- A quantitative evaluation of category-level object grasping with VS used for end-effector positioning. Grasping unseen instances of the categories shoe, toy car, and mug with a mean success rate of 90% ( $n = 30$ ).

The remainder of this paper is organized as follows: Section II reviews related work, Section III presents our novel ViT-VS methodology, Section IV details our experimental evaluation and results, and Section V concludes the paper with a discussion of our findings.

## II. RELATED WORK

This section introduces classical visual servoing schemes, Section II-A, Section II-B reviews emerging deep learning approaches, and Section II-C discusses the potential of ViT features to advance visual servoing.

### A. Classical Visual Servoing

Classical visual servoing can be broadly classified into IBVS [19] and PBVS [20]. These approaches rely on hand-crafted geometric features such as points, lines, and moments for robot guidance.

PBVS requires both camera intrinsic parameters and the object’s 3D model, and while it can theoretically achieve global asymptotic stability, it relies heavily on accurate pose estimation. In contrast, IBVS operates directly in the image space, requiring only camera intrinsic parameters. IBVS demonstrates robust performance against calibration errors and image noise, though it can only guarantee local asymptotic stability [1], [2]. Traditional feature detectors have been utilized for IBVS, with SIFT features performing end-effector positioning [7] and SURF features executing object-specific grasping [6]. However, while these traditional feature extracting methods offer general applicability, they have been observed to struggle with occlusions, varying illumination, and complex environments [21]. In order to avoid explicit feature extraction Direct Visual Servoing (DVS) was introduced, while achieving lower positioning error compared to classical approaches, it suffers from a limited convergence domain [8]. These limitations in handling real-world complexities for classical VS strategies have motivated the exploration of learning-based approaches.

### B. Learning Based Visual Servoing

In recent years multiple works contributed to visual servoing by utilizing deep learning, overcoming the problematic of occlusion, lighting variations, scene changes, and image perturbations. The authors in [9] combine classical PBVS with convolutional neural networks to regress camera pose using synthetically generated datasets for training, demonstrating robust convergence against environmental variations. Building on pose-based approaches, [12] explores deep metric learning by creating a common latent space for camera poses and image representations, incorporating perturbed

samples in the training data for enhanced robustness. Several works utilize siamese networks for visual servoing, in [10] directly regressing camera velocity without pose estimation, while the method of [14] jointly learns feature extraction and transformation through 3D equivariance constraints for wide-baseline visual servoing. Alternative feature-based approaches include the work of [13], which uses an unsupervised convolutional autoencoder to learn compact image representations that generalize to similar unseen targets. In the research of [11] classical IBVS is combined with neural networks for feature extraction and matching, reaching promising final positioning error though requiring a rendering engine and object model. Similarly, [5] employs a simulation-to-real transfer approach using object models, learning end-to-end robotic motion and enabling direct deployment on real robots after training purely in simulation. A recent work [15] replaces pixel brightness with neural network feature representations for DVS, though still facing challenges with illumination variations. While these learning-based approaches demonstrate significant improvements over classical methods, they either require task-specific training and data generation [9], [10], [13], [14], [12] or target object models [11], [5], limiting their immediate practical deployment. This highlights an ongoing challenge in developing more flexible and readily deployable visual servoing solutions.

### C. Foundation Model Features

Vision Transformers have recently emerged as powerful architectures for visual tasks, offering robust semantic understanding through self-attention mechanisms [16]. Pre-trained ViTs demonstrate strong zero-shot capabilities and generalization across related object categories [22], [23]. While primarily used for classification and detection tasks, their ability to extract general features makes them promising for visual servoing applications.

In this work, we leverage the suitability of ViT features for zero-shot vision tasks. By using ViT features for IBVS we combine the advantages of classical and learning-based VS. The method presented in the next section combines the advantages of classical approaches, i.e. general applicability but no need for offline training or fine tuning, with the advantages of deep learning approaches, i.e. robustness to occlusion, lighting variations and image perturbations.

## III. DEEP GENERALIZABLE VISUAL SERVOING

Fig. 2 illustrates our deep zero-shot visual servoing pipeline, which combines ViT correspondence matching, initial rotation compensation, IBVS control and velocity stabilization. Our approach leverages pretrained DINOv2 [17] models to extract patch embeddings. Correspondences are established using cosine similarity and cyclical distance metric. To maintain spatial diversity across the image the correspondences are randomly selected from top-K matches. To handle the inherent rotation invariance of ViTs, we implement an initial rotation compensation to evaluate the best initial pose. The estimated patch correspondences are then used as input

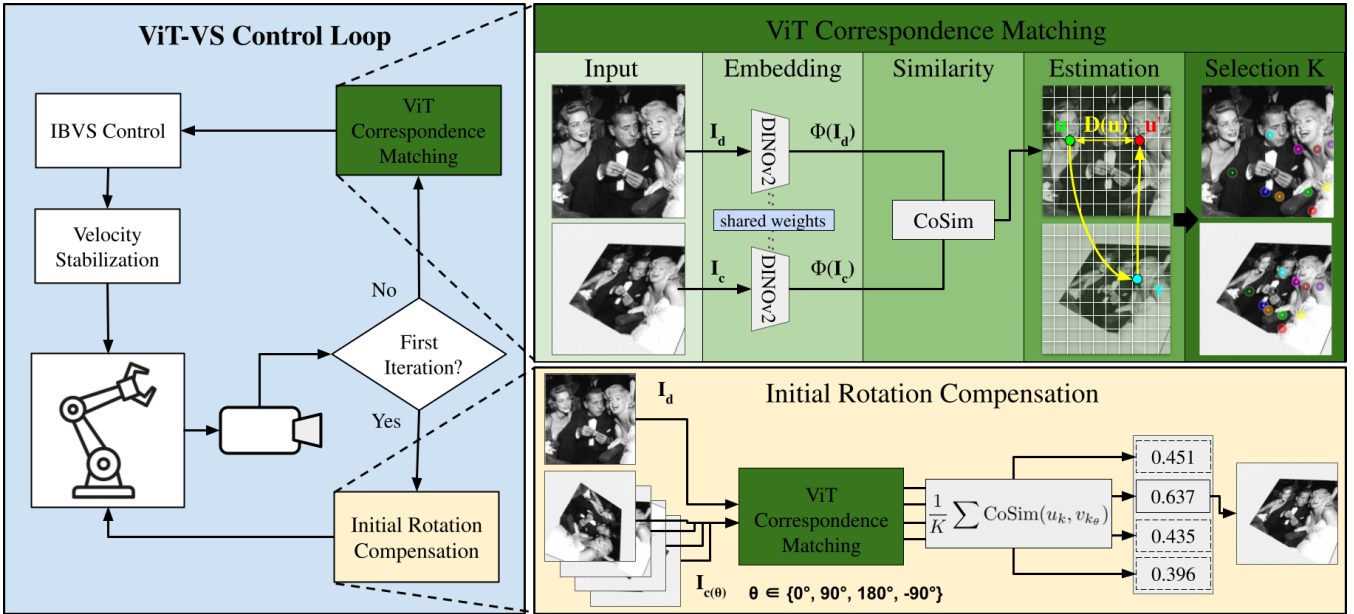


Fig. 2. **ViT-VS Overview** Left: Visual servoing control loop integrating initial rotation compensation and IBVS control with ViT correspondences and velocity stabilization. Top right: Feature correspondence pipeline using DINOv2-based patch embeddings from desired and current image, where matches are estimated through cosine similarity and cyclical distance metrics, with random selection from top-K matches for spatial diversity. Bottom right: Initial rotation compensation mechanism evaluating four discrete rotations ( $0^\circ$ ,  $90^\circ$ ,  $180^\circ$ ,  $-90^\circ$ ) to determine optimal starting pose through mean feature similarity.

to the classical IBVS controller, with additional velocity stabilization through exponential moving average filtering to ensure smooth trajectories. In the following sections, we describe each component of our pipeline in detail.

### A. ViT Correspondence Matching

As visualized in the top right block of Fig. 2 our method takes a desired image  $I_d$  and current image  $I_c$  as input. Using the DINOv2 small model [17], we extract descriptor sets  $\Phi(I_d), \Phi(I_c) \in \mathbb{R}^{H' \times W' \times D}$ . Following [23] and [24], we adopt best buddy pairs matching concept from [25], which provides robust feature correspondences between images. Starting with a point  $u \in \{1, \dots, H'\} \times \{1, \dots, W'\}$  in  $\Phi(I_d)$ , we find its nearest neighbor  $v$  in  $\Phi(I_c)$  using cosine similarity:

$$v = \arg \max_w \text{CoSim}(\Phi(I_d)_u, \Phi(I_c)_w) \quad (1)$$

To verify this match, we find the nearest neighbor  $u'$  of  $v$  back in  $I_d$ :

$$u' = \arg \max_z \text{CoSim}(\Phi(I_d)_z, \Phi(I_c)_v) \quad (2)$$

Using these matched pairs, we construct a cyclical distance map  $D \in \mathbb{R}^{H' \times W'}$  as:

$$D_u = -\|u - u'\|_2 \quad (3)$$

This cyclical distance as given in Equation 3 improves upon [25] by enabling the consistent selection of  $K$  semantic correspondences while incorporating spatial priors. A cyclical distance of zero indicates a perfect match where the correspondence maps back to the original point. We randomly select  $K$  correspondences that exceed our cyclical distance threshold, promoting better spatial distribution of

features across the image. This distribution strategy proves to be crucial for robust convergence by preventing feature concentration in visually distinctive regions.

**Feature Binning:** To enhance feature robustness, we implement hierarchical feature binning as introduced by [23]. The binning hierarchy parameter  $\beta$  determines the contextual scope around each patch. At  $\beta = 1$ , each patch is combined with its eight immediate neighbors in a  $3 \times 3$  grid. Each increment of  $\beta$  adds a new ring of context, with average pooling ensuring smooth feature transitions. This hierarchical approach enriches the descriptors with broader contextual information at the cost of computational complexity.

**Operating Resolution:** DINOv2 operates at input resolutions between  $224 \times 224$  and  $518 \times 518$  pixels. Due to the computational demands of ViTs, feature extraction is limited to this resolution range. Input images must be resized within this range, directly affecting both the granularity of extractable features and real-time performance capabilities.

**Foreground Segmentation** In cases where the desired image includes background, it is necessary to utilize a foreground segmentation for the creation of segmentation masks. For this purpose we utilize Segment Anything [26].

### B. Initial Rotation Compensation

Vision Transformers, including DINOv2 [17], are not inherently rotation invariant due to their architecture with fixed position encodings and patch-based processing [16]. To address this limitation, we implement a rotation compensation step as shown in the bottom right of Fig. 2.

We evaluate the correspondence matching at four rotations ( $0^\circ$ ,  $90^\circ$ ,  $180^\circ$ ,  $-90^\circ$ ), calculating the mean cosine similarity score between the selected feature pairs  $(u_k, v_k)$ . As shown in Equation 4, we determine the optimal initial rotation angle

$\theta^*$  that maximizes the mean similarity score in all matched pairs.

$$\theta^* = \underset{\theta \in \{0^\circ, 90^\circ, 180^\circ, -90^\circ\}}{\arg \max} \left( \frac{1}{K} \sum_{k=1}^K \text{CoSim}(u_k, v_k^\theta) \right) \quad (4)$$

This optimal rotation is applied to the robot, aligning the current image orientation with the desired image before initiating the visual servoing process.

### C. Image Based Visual Servoing with ViT

Following rotation compensation, we implement classical IBVS as described by [1]. The control law aims to minimize the error:

$$e(t) = s(m(t), a) - s^* \quad (5)$$

where  $s(m(t), a)$  represents the current image feature vector extracted from the image measurements  $m(t)$  using camera intrinsic parameters  $a$ , and  $s^*$  denotes the desired features. For  $n$  feature points obtained through ViT matching, we compute the velocity control law as given in Equation 6.

$$v_c = -\lambda \hat{L}_e^+ e \quad (6)$$

Here,  $\hat{L}_e^+$  is the Moore-Penrose pseudoinverse of the interaction matrix, which we approximate using the depth values  $Z_i$  of the corresponding feature points from the current depth-image. The parameter  $\lambda$  is included to ensure exponential error decrease. For each feature point, the interaction matrix  $L_i$  related to  $s_i = (x_i, y_i)$  is constructed as:

$$L_i = \begin{bmatrix} -\frac{1}{Z_i} & 0 & \frac{x_i}{Z_i} & x_i y_i & -(1 + x_i^2) & y_i \\ 0 & -\frac{1}{Z_i} & \frac{y_i}{Z_i} & 1 + y_i^2 & -x_i y_i & -x_i \end{bmatrix} \quad (7)$$

where  $(x, y)$  are feature point coordinates transformed to real-world units using camera intrinsics. The complete interaction matrix combines these individual matrices as:

$$L_e = \begin{bmatrix} L_1 \\ \vdots \\ L_{n-1} \\ L_n \end{bmatrix} \quad (8)$$

The control process follows three steps: first, calculating the error  $e(t)$  as defined in Equation 5, then computing the interaction matrix  $L_e$  as given in Equations 7 and 8, and finally determining the velocity vector  $v_c$  in Equation 6. The resulting velocity vector contains six components describing the cameras' linear and angular velocities in the camera coordinate frame.

### D. Velocity Stabilization

Visual servoing systems face multiple uncertainty sources: feature detection variations, depth estimation inaccuracies, and numerical instabilities in matrix computations. While our randomized feature selection improves convergence reliability, it can introduce velocity fluctuations affecting trajectory smoothness and mechanical stability. To address this, we

implement an Exponential Moving Average (EMA) filter for each velocity component:

$$v_t = \alpha v_{\text{new}} + (1 - \alpha)v_{t-1} \quad (9)$$

where  $v_t$  is the smoothed velocity,  $v_{\text{new}}$  is the newly computed velocity,  $v_{t-1}$  is the previous smoothed velocity, and  $\alpha$  is the smoothing factor. This filtering approach effectively dampens unwanted velocity fluctuations while maintaining system responsiveness, ensuring smooth trajectories throughout the servoing process.

## IV. EXPERIMENTS

Our experimental evaluation consists of simulation studies, detailed system analysis, and extensive robotic experiments across different application scenarios. We first describe our implementation details, followed by simulation experiments that benchmark our method against state-of-the-art approaches. We then present system analysis results and conclude with real-world robotic experiments in three different scenarios.

### A. Experimental Setup

**Hardware Setup** All experiments are conducted using a consistent hardware and software setup across simulation and real-world scenarios. The simulation experiments are conducted on a NVIDIA RTX 4060 Ti 16 GB and Intel Core i7-13700K CPU. Robotic experiments are conducted on an NVIDIA RTX 4070 mobile GPU and AMD Ryzen 9 7940HS CPU. The mobile manipulator platform consists of a Universal Robots UR5 mounted on a Mobile Industrial Robots MiR100. Industrial manipulation experiments are performed with a custom gripper and object grasping is performed with a Robotiq 2F-85. An Intel RealSense D435i camera is used for all experiments.

**ViT Model Configuration** Our experimental setup employs DINOv2 [17], pretrained on ImageNet1k [27], ViT-Small/14 architecture for feature extraction, using layer 11 for token-based features. The model operates with a patch size and stride of 14 pixels. Feature binning is adopted from [23]. Depending on the experiment, we use  $\beta = 1$  or  $\beta = 2$  binning hierarchies combined with a DINOv2 input resolution between  $224 \times 224$  and  $308 \times 308$  pixels as described in Section III-A. However, if not stated otherwise, we use  $308 \times 308$  pixels and  $\beta = 1$ . For each iteration, we extract 24 feature pairs for matching between current and desired view using the approach introduced in Section III-A. For similarity estimation and cyclic distance map computation we use the implementation of [24], and for final correspondence selection we utilize a cutoff threshold of 1.

**Simulation Environment** The simulation environment is replicated from Deep Metric Learning for Visual Servoing (DMLVS) [12]. A virtual Intel RealSense D435i camera with a resolution of  $640 \times 480$  pixels is used, the target is the ‘‘hollywood poster’’ in  $60 \times 80\text{cm}$ , visible in Fig. 3. We closely follow the approach of [12] and generate 500 distinct initial camera poses. Camera positions are sampled within a cuboid of  $1.2\text{m} \times 1.2\text{m} \times 0.3\text{m}$  volume centered on

the desired position. The look-at points are distributed across four concentric circles on the poster plane, with radii of 8, 16, 24, and 32cm from the poster’s center. Camera orientations are samples using the look-at function, with an additional random rotation around the focal axis within  $[-120^\circ, 120^\circ]$ . The desired camera position is set at 0.6m elevation from the poster center. This configuration yields average initial position errors of  $46.42 \pm 16.99cm$  and orientation errors of  $74.12 \pm 27.71^\circ$ .

**Perturbation Settings** To evaluate the robustness of our method and compare it to the deep learning results presented in [12] we conduct our experiments with unperturbed and perturbed images. The perturbation parameters are taken from the codebase of [12], and are implemented using Torchvision transforms: Colorjitter with a brightness of 0.6 and contrast of 0.4. The random erasing with a probability of 0.5 on a scale of 0.02 to 0.33 and a ratio of 0.3 to 3.3. The Gaussian blur is implemented with a mean of 0 and a sigma of 0.05.

**Evaluation Metrics** Convergence is reached when the velocities are close to zero, and initial position and rotation error are reduced by more than 90%, as done in [12]. Furthermore we also report the Absolute Pose Error (APE) and length ratio; APE quantifies the cumulated error in relation to the optimal PBVS trajectory and the length ratio quantifies the ratio between executed and ideal trajectory.

### B. Simulation Experiments

This section compares ViT-VS with both classical and learning-based visual servoing approaches, in simulation. To ensure fair comparison with classical ones, ViT-VS’ matching strategy is replicated. Feature matching with SIFT [28], a floating-point descriptors, is done using the Euclidean distance, and for ORB [29] and AKAZE [30], both binary descriptors, the Hamming distance is used. Matches are ranked by descriptor distance, and 24 matches from the top-ranked candidates, as is the case for ViT-VS, are used for servoing. Table I shows that ViT-VS converges with a success rate of 100% in the unperturbed scenarios. Hence, the presented convergence rate is on par with the state-of-the-art learning-based DMLVS [12], but without object- or scene-specific finetuning, and also significantly higher than that of classical descriptors. Under image perturbations as exemplified in Fig. 3, ViT-VS achieves a 76.6% success rate, improving over all classical feature-based methods and deep learning-based methods. The end error of ViT-VS, despite being inferior to classical methods, is competitive to PBVS approaches, in the case of image perturbations. The translational end error is close to that of DMLVS, with  $19.29 \pm 12.81cm$  for DMLVS and  $21.54 \pm 12.11cm$  for ViT-VS, while the rotational end error is better, with  $1.92 \pm 1.28^\circ$  for DMLVS compared to  $1.83 \pm 0.98^\circ$  for ViT-VS. The high end error in comparison to classical approaches is a consequence of the coarse feature maps of ViTs, which is 1/14 of the input resolution in the case of our configuration with DINOv2-small. This leads to correspondences being matched in a space with a resolution of  $22 \times 22$ , as compared to classical IBVS



Fig. 3. ViT feature correspondences between desired image (left) and current perturbed image (right) in simulation environment.

methods matching in a space with the resolution of the input image. Deep learned pose-based methods lead to a better APE and trajectory length since these metrics are designed for evaluating PBVS behavior. ViT-VS achieves translational APE comparable to the best performing classical descriptor, ORB,  $17.14 \pm 6.65cm$  compared to  $16.60 \pm 5.66cm$ , the best rotational accuracy  $16.34 \pm 5.05^\circ$  among non-finetuned methods, and the best length ratio at  $1.21 \pm 0.39$ .

Hence, ViT-VS combines the advantages of classical IBVS approaches, universal and general applicability since not requiring finetuning, and those of PBVS, high convergence rates, and better APE and length ratios than classical approaches. This makes our method ideal for real-world robotic manipulation tasks where generality is required, and consistent convergence outweighs sub-millimeter accuracy.

**Initial Rotation Compensation** Table I shows that the rotation compensation improves ViT-VS performance significantly. With this mechanism ViT-VS achieves 100.0% and 76.6%, and without 83.8% and 57.2% convergence rates are achieved, for unperturbed and perturbed images, respectively. The compensation allows our method to achieve state-of-the-art convergence rates while operating without object or scene-specific finetuning.

**Frame Rate Analysis** In Fig. 4 we presents the frame rate for ViT-VS, averaged over 100 runs using different configurations of DINOv2 [17]. Our experiments show that feature binning as introduced in Section III-A has a larger impact on computational efficiency than the choice of DINOv2 backbone size. Based on these results, we focus on two configurations: DINOv2-Small with  $224 \times 224$  pixel input and  $\beta = 2$  binning, or  $308 \times 308$  pixel input with  $\beta = 1$  binning. While the presented frame rates leave room for improvement, our robotic experiments show that the application of trajectory regularization effectively addresses motion jitter, resulting in smoother trajectories, as detailed in the following subsection.

**Trajectory Regularization** The influence of the trajectory smoothing parameter  $\alpha$  as defined in III-D, in a range from 0.5 to 0.9, is evaluated and illustrated in Fig. 5. Lower  $\alpha$  values lead to reduced length ratios, however, they also lead to an increase in end-positioning error. We choose an  $\alpha$  of 0.8 as the standard configuration. This value choice balances length ratio and end error.

TABLE I  
COMPARISON TO THE STATE OF THE ART RESULTS MARKED WITH AN ASTERISK ARE TAKEN FROM [12].

	Method	Perturbed image	Converged [%]	End error [mm]	End error [°]	APE [cm]	APE [°]	Length ratio
Deep finetuning	AEVS [13]*	×	33.6	0.01 ± 0.00	0.00 ± 0.00	2.74 ± 6.36	2.66 ± 6.4	2.75 ± 4.85
	PBVS-CNN, e.g., [9]*	×	75.6	33.52 ± 6.45	1.71 ± 0.65	3.13 ± 1.04	1.85 ± 0.96	1.11 ± 0.08
	PBVS-CNN, e.g., [9]*	✓	36.8	32.21 ± 15.71	2.37 ± 1.57	4.00 ± 0.74	<b>2.55 ± 0.64</b>	<b>1.12 ± 0.10</b>
	DMLVS, $K = 50$ , [12]*	×	100.0	0.04 ± 0.03	0.00 ± 0.00	4.00 ± 0.72	1.08 ± 2.50	1.18 ± 0.23
	DMLVS, $K = 50$ , [12]*	✓	<b>76.0</b>	<b>19.29 ± 12.81</b>	<b>1.92 ± 1.28</b>	<b>3.31 ± 0.61</b>	3.72 ± 1.60	1.14 ± 0.11
No finetuning	DVS [8]*	×	9.8	0.00 ± 0.00	0.00 ± 0.00	17.32 ± 12.48	31.0 ± 14.35	2.60 ± 4.90
	SIFT IBVS	×	89.6	1.17 ± 2.33	0.09 ± 0.11	18.22 ± 7.25	27.76 ± 11.31	2.37 ± 1.74
	SIFT IBVS	✓	24.0	2.78 ± 5.32	0.23 ± 0.47	20.29 ± 9.22	26.66 ± 11.15	3.64 ± 2.87
	ORB IBVS	×	98.6	3.32 ± 1.49	0.25 ± 0.12	16.60 ± 5.66	25.66 ± 10.64	1.82 ± 1.19
	ORB IBVS	✓	58.4	3.86 ± 3.36	0.30 ± 0.26	<b>16.76 ± 6.14</b>	24.33 ± 10.46	1.91 ± 1.18
	AKAZE IBVS	×	89.0	1.03 ± 1.21	0.08 ± 0.11	16.99 ± 5.90	28.79 ± 11.86	1.91 ± 1.26
	AKAZE IBVS	✓	58.0	<b>1.44 ± 1.06</b>	<b>0.12 ± 0.09</b>	16.86 ± 5.87	25.70 ± 10.71	1.91 ± 1.61
	ViT-VS (no rot. comp.)	×	83.8	21.68 ± 9.20	1.66 ± 0.69	18.36 ± 6.65	25.37 ± 9.61	2.34 ± 2.98
	ViT-VS (no rot. comp.)	✓	57.2	24.10 ± 11.80	1.94 ± 1.01	17.94 ± 6.38	22.99 ± 8.76	2.95 ± 3.40
	<b>ViT-VS (ours)</b>	×	<b>100.0</b>	<b>18.62 ± 10.69</b>	<b>1.50 ± 0.78</b>	<b>17.14 ± 6.65</b>	<b>16.34 ± 5.05</b>	<b>1.21 ± 0.39</b>
	<b>ViT-VS (ours)</b>	✓	<b>76.6</b>	21.54 ± 12.11	1.83 ± 0.98	17.05 ± 6.18	<b>16.29 ± 5.24</b>	<b>1.90 ± 1.44</b>

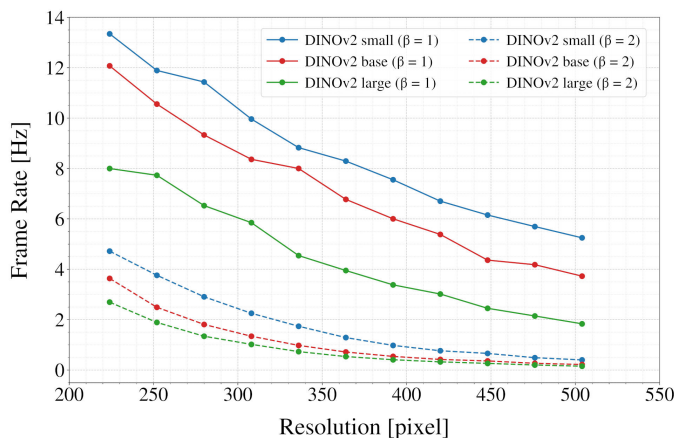


Fig. 4. ViT-VS Frame Rate Analysis using NVIDIA RTX 4070 Mobile GPU across various model configurations.

### C. Robotic Experiments

This section shows a detailed robotic experiment on the image that was used for comparing ViT-VS’ convergence to the state of the art in simulation. Following that, experiments are provided that show that the presented method is suited for robustly compensating a mobile robot’s positioning error for box manipulation. Ultimately, we demonstrate that ViTs are semantically stable in a way to enable category-level object grasping.

**Detailed Robotic Experiment** This experiment demonstrates real-world evaluation on the “hollywood poster” for 1500 iterations. The initial position error is  $\Delta \mathbf{r}_0 = (-45.60\text{cm}, 18.63\text{cm}, -11.21\text{cm}, 10.17^\circ, -15.48^\circ, -153.08^\circ)$ . Fig. 6 visualizes the initial image (a), the desired image (b) and the final image (c). Fig. 6(d) and (e) present the camera velocities and position and rotation errors. The initial rotation compensation is not indicated since it is not part of the control loop. The best rotation was found to be a  $180^\circ$ . The final position error is  $\Delta \mathbf{r}_{\text{final}} = (0.38\text{cm}, 0.44\text{cm}, -0.25\text{cm}, 0.44^\circ, -0.54^\circ, -0.40^\circ)$ ,

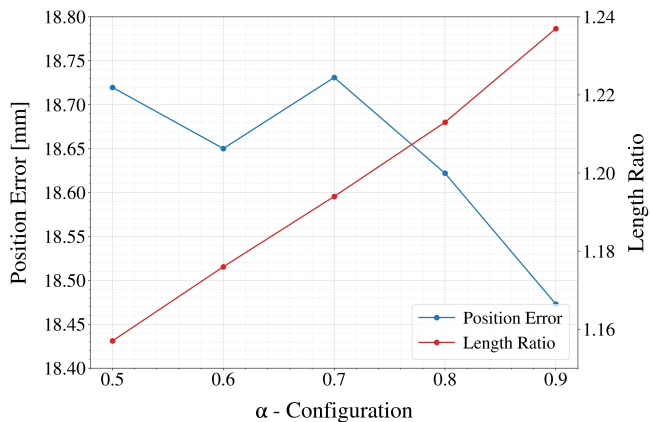


Fig. 5. Alpha evaluation based on simulation runs without perturbation.

showcasing full convergence from an only partially visible and heavily rotated initial position.

**Industrial Use-case** This experiment demonstrates the suitability of our approach for an industrial use case with a mobile manipulator. The mobile robot navigates to a work cell using ROS navigation stack [31] and two laser scanners. This results in a positioning uncertainty of  $\pm 10\text{cm}$  of the MiR100 base. Fig. 7 shows the desired image (a), an initial image (b), and the mobile robot after convergence (c). The trials are performed using boxes with different appearances, e.g., boxes with missing labels or structural differences. We achieve a 100% success rate on 20 trials of box lifting. The convergence behavior and positioning errors for all trials are visualized in Fig. 7(d). For this experiment  $\beta = 2$  and an image resolution of  $224 \times 224$  pixels is used with DINOv2-Small to focus on geometry rather than texture.

**Category-level Object Grasping** Category-level object sorting experiments showcase the strong generalization capability of our method, allowing to perform pick and place tasks. Table II reports the success rates for grasping and sorting of randomly placed singulated objects. The left part of Fig. 8 shows the reference objects and the grasp points of

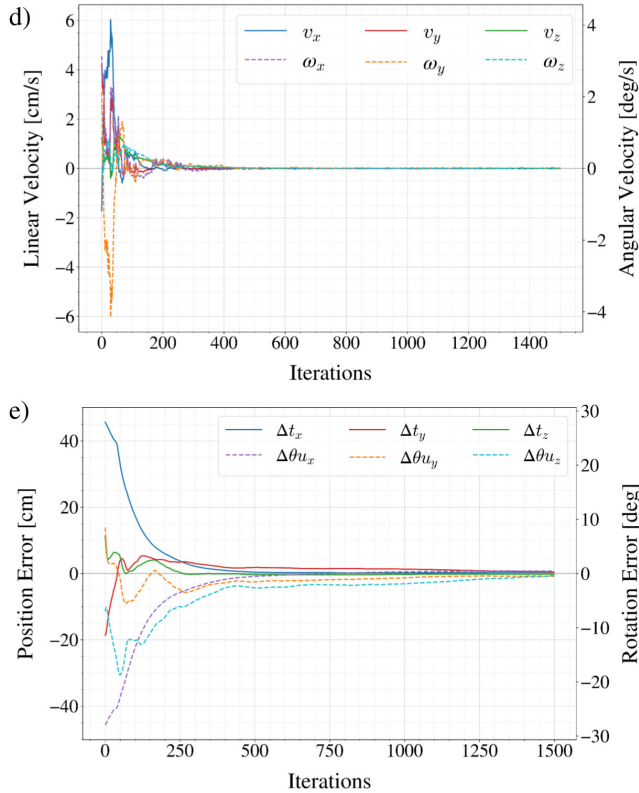
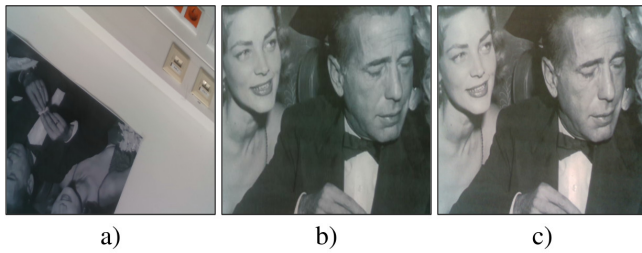


Fig. 6. **Detailed robotic experiment** a) Initial image. b) Desired image. c) Final image. d) Camera velocities. e) Pose difference.

TABLE II  
CATEGORY-LEVEL OBJECT GRASPING EXPERIMENTS SUCCESSES FOR GRASPING OF SINGULATED UNSEEN OBJECT INSTANCES OF THE CATEGORIES.

Object	Car		Shoe		Mug	
	blue	black	blue	green	black	white
Successes	5/5	3/5	5/5	5/5	4/5	5/5

the corresponding objects, the right part shows the unseen instance of the three categories mug, toy car, and shoe. The left column shows the reference objects and the grasp points of the corresponding objects. For each unseen object instance 5 picking tries are performed; resulting in 10 tries per object category. The initial position is chosen to capture the full table plan. Objects are separated from the background using Segment Anything [26]. ViT-VS, with a model configuration of  $\beta = 2$  and  $224 \times 224$  input resolution, positions the end effector relative to the object for triggering the grasping motion, which is predefined for the seen object instance.

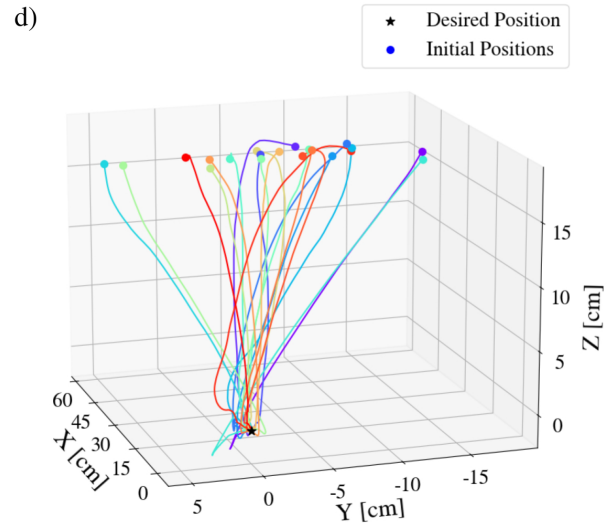
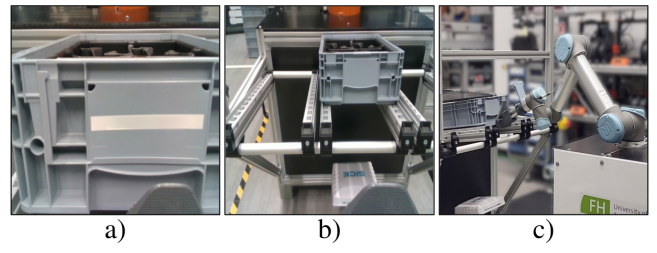


Fig. 7. **Industrial use-case** a) Desired image. b) Example for initial image. c) External view before gripping. d) Tool center point trajectories plot aligned at goal position.

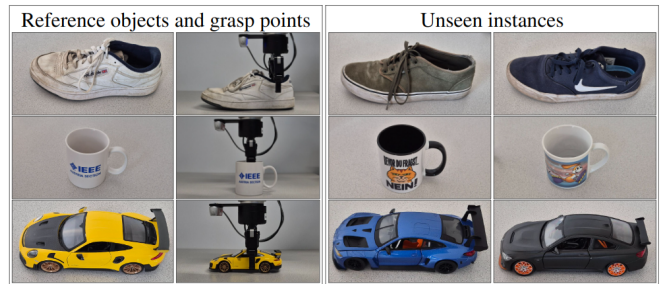


Fig. 8. **Objects for sorting** Left side: Reference objects with corresponding grasp points. Right side: unseen object instances.

A successful grasp requires the robot to grasp and lift the object. Our method demonstrates robust performance over all object categories, achieving success rates of 100% for shoes, 90% for mugs, and 80% for toy cars. The failed attempt for the mug occurred because the mug slipped out of the gripper after successful convergence and grasp. The two failed attempts of the toy cars occurred due to table plane collisions; one due to bad convergence, one due to bad convergence caused by the initial rotation compensation retrieving the incorrect rotation. Fig. 9(a)-(e) shows a representative grasping sequence of the blue unseen toy car.

## V. CONCLUSIONS

This work demonstrates the advantages of pretrained Vision Transformer features for visual servoing; Convergence

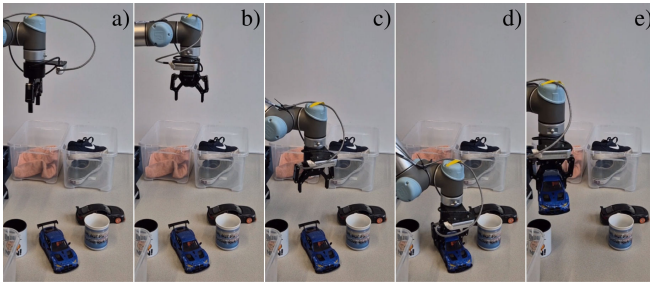


Fig. 9. **Unseen object manipulation** Grasping of an unseen instance of the category toy car: a) initial position, b) compensated rotation, c) converged position, d) grasping, and e) object manipulation.

rates are comparable to learning-based methods, yet features are generally applicable without finetuning, as is the case for classical image-based visual servoing. Diverse robotics experiments demonstrate the usefulness and generality of Vision Transformer features for industrial tasks and household tasks, such as object manipulation and unseen object instance grasping. Future work will investigate strategies for improving the positioning errors which are dictated by the resolution of the Vision Transformers' feature maps.

## REFERENCES

- [1] F. Chaumette and S. Hutchinson, "Visual servo control. i. basic approaches," *IEEE Robotics & Automation Magazine*, vol. 13, no. 4, pp. 82–90, 2006.
- [2] C. Francois and H. Seth, "Visual servo control part ii: Advanced approaches," *IEEE Trans on Robotics and Automation*, vol. 14, no. 1, pp. 109–118, 2007.
- [3] Y. Chen, Y. Wu, Z. Zhang, Z. Miao, H. Zhong, H. Zhang, and Y. Wang, "Image-based visual servoing of unmanned aerial manipulators for tracking and grasping a moving target," *IEEE Transactions on Industrial Informatics*, vol. 19, no. 8, pp. 8889–8899, 2022.
- [4] C. De Farias, M. Adjigle, B. Tamadazte, R. Stolkin, and N. Marturi, "Dual quaternion-based visual servoing for grasping moving objects," in *2021 IEEE 17th International Conference on Automation Science and Engineering (CASE)*. IEEE, 2021, pp. 151–158.
- [5] E. Y. Puang, K. P. Tee, and W. Jing, "Kovis: Keypoint-based visual servoing with zero-shot sim-to-real transfer for robotics manipulation," in *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2020, pp. 7527–7533.
- [6] T. La Anh and J.-B. Song, "Robotic grasping based on efficient tracking and visual servoing using local feature descriptors," *International Journal of Precision Engineering and Manufacturing*, vol. 13, pp. 387–393, 2012.
- [7] F. Hoffmann, T. Nierobisch, T. Seyffarth, and G. Rudolph, "Visual servoing with moments of sift features," in *2006 IEEE International Conference on Systems, Man and Cybernetics*, vol. 5. IEEE, 2006, pp. 4262–4267.
- [8] C. Collewet, E. Marchand, and F. Chaumette, "Visual servoing set free from image processing," in *2008 IEEE International Conference on Robotics and Automation*, 2008, pp. 81–86.
- [9] Q. Bateux, E. Marchand, J. Leitner, F. Chaumette, and P. Corke, "Training deep neural networks for visual servoing," in *2018 IEEE International Conference on Robotics and Automation (ICRA)*, 2018, pp. 3307–3314.
- [10] S. Felton, E. Fromont, and E. Marchand, "Siame-se(3): regression in se(3) for end-to-end visual servoing," in *2021 IEEE International Conference on Robotics and Automation (ICRA)*, 2021, pp. 14454–14460.
- [11] N. Adrian, V.-T. Do, and Q.-C. Pham, "Dfbvs: Deep feature-based visual servo," in *2022 IEEE 18th International Conference on Automation Science and Engineering (CASE)*, 2022, pp. 1783–1789.
- [12] S. Felton, E. Fromont, and E. Marchand, "Deep metric learning for visual servoing: when pose and image meet in latent space," in *2023 IEEE International Conference on Robotics and Automation (ICRA)*, 2023, pp. 741–747.
- [13] S. Felton, P. Brault, E. Fromont, and E. Marchand, "Visual servoing in autoencoder latent space," *IEEE Robotics and Automation Letters*, vol. 7, no. 2, pp. 3234–3241, 2022.
- [14] J. Huh, J. Hong, S. Garg, H. S. Park, and V. Isler, "Self-supervised wide baseline visual servoing via 3d equivariance," in *2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2022, pp. 2227–2233.
- [15] M. Quaccia, A. N. André, Y. Yoshiyasu, and G. Caron, "A study on learned feature maps toward direct visual servoing," in *2024 IEEE/SICE International Symposium on System Integration (SII)*, 2024, pp. 520–525.
- [16] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby, "An image is worth 16x16 words: Transformers for image recognition at scale," in *International Conference on Learning Representations*, 2021. [Online]. Available: <https://openreview.net/forum?id=YicbFdNTTy>
- [17] M. Oquab, T. Darcet, T. Moutakanni, H. V. Vo, M. Szafraniec, V. Khalidov, P. Fernandez, D. HAZIZA, F. Massa, A. El-Nouby, M. Assran, N. Ballas, W. Galuba, R. Howes, P.-Y. Huang, S.-W. Li, I. Misra, M. Rabbat, V. Sharma, G. Synnaeve, H. Xu, H. Jegou, J. Mairal, P. Labatut, A. Joulin, and P. Bojanowski, "DINOv2: Learning robust visual features without supervision," *Transactions on Machine Learning Research*, 2024, featured Certification. [Online]. Available: <https://openreview.net/forum?id=a68SUt6zFt>
- [18] J. Zhang, C. Herrmann, J. Hur, E. Chen, V. Jampani, D. Sun, and M.-H. Yang, "Telling left from right: Identifying geometry-aware semantic correspondence," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 3076–3085.
- [19] S. Hutchinson, G. D. Hager, and P. I. Corke, "A tutorial on visual servo control," *IEEE transactions on robotics and automation*, vol. 12, no. 5, pp. 651–670, 1996.
- [20] W. J. Wilson, C. W. Hulls, and G. S. Bell, "Relative end-effector control using cartesian position based visual servoing," *IEEE Transactions on Robotics and Automation*, vol. 12, no. 5, pp. 684–696, 1996.
- [21] E. Karami, S. Prasad, and M. Shehata, "Image matching using sift, surf, brief and orb: performance comparison for distorted images," *arXiv preprint arXiv:1710.02726*, 2017.
- [22] M. Caron, H. Touvron, I. Misra, H. Jegou, J. Mairal, P. Bojanowski, and A. Joulin, "Emerging properties in self-supervised vision transformers," in *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021, pp. 9630–9640.
- [23] S. Amir, Y. Gandelsman, S. Bagon, and T. Dekel, "Deep vit features as dense visual descriptors," *ECCVW What is Motion For?*, 2022.
- [24] W. Goodwin, S. Vaze, I. Havoutis, and I. Posner, "Zero-shot category-level object pose estimation," in *European Conference on Computer Vision*. Springer, 2022, pp. 516–532.
- [25] T. Dekel, S. Oron, M. Rubinstein, S. Avidan, and W. T. Freeman, "Best-buddies similarity for robust template matching," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 2021–2029.
- [26] A. Kirillov, E. Mintun, N. Ravi, H. Mao, C. Rolland, L. Gustafson, T. Xiao, S. Whitehead, A. C. Berg, W.-Y. Lo *et al.*, "Segment anything," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 4015–4026.
- [27] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," *Advances in neural information processing systems*, vol. 25, 2012.
- [28] D. G. Lowe, "Distinctive image features from scale-invariant keypoints," *International journal of computer vision*, vol. 60, pp. 91–110, 2004.
- [29] E. Rublee, V. Rabaud, K. Konolige, and G. Bradski, "Orb: An efficient alternative to sift or surf," in *2011 International conference on computer vision*. Ieee, 2011, pp. 2564–2571.
- [30] P. F. Alcantarilla and T. Solutions, "Fast explicit diffusion for accelerated features in nonlinear scale spaces," *IEEE Trans. Patt. Anal. Mach. Intell.*, vol. 34, no. 7, pp. 1281–1298, 2011.
- [31] E. Marder-Eppstein, E. Berger, T. Foote, B. Gerkey, and K. Konolige, "The office marathon: Robust navigation in an indoor office environment," in *2010 IEEE international conference on robotics and automation*. IEEE, 2010, pp. 300–307.