

Push to know! - Visuo-Tactile based Active Object Parameter Inference with Dual Differentiable Filtering

Anirvan Dutta, Etienne Burdet and Mohsen Kaboli

Abstract—For robotic systems to interact with objects in dynamic environments, it is essential to perceive the physical properties of the objects such as shape, friction coefficient, mass, center of mass, and inertia. This not only eases selecting manipulation action but also ensures the task is performed as desired. However, estimating the physical properties of especially novel objects is a challenging problem, using either vision or tactile sensing. In this work, we propose a novel framework to estimate key object parameters using non-prehensile manipulation using vision and tactile sensing. Our proposed active dual differentiable filtering (ADDF) approach as part of our framework learns the object-robot interaction during non-prehensile object push to infer the object’s parameters. Our proposed method enables the robotic system to employ vision and tactile information to interactively explore a novel object via non-prehensile object push. The novel proposed N -step active formulation within the differentiable filtering facilitates efficient learning of the object-robot interaction model and during inference by selecting the next best exploratory push actions (where to push? and how to push?). We extensively evaluated our framework in simulation and real-robotic scenarios, yielding superior performance to the state-of-the-art baseline.

I. INTRODUCTION

To manipulate novel objects, humans often perceive object properties through actions such as holding, grasping or pushing to gain better control [1], [2]. In such *interactive visuo-tactile perception*, active physical interactions or explorations are made to enhance object perception [3]–[6]. In this work, we investigate inferring physical object properties (such as friction coefficient, mass, center of mass, and inertia) with a robotic system using non-prehensile push actions via visuo-tactile sensory information (see Fig. 1).

Pushing an object to explore its parameters is simpler than grasping or lifting [7], especially when the objects are large and heavy or when no prior knowledge about the object is present. Further, in grasping or lifting, object geometry, grasp stability and other factors come into play which is not the case for pushing. Moreover, since the object is not rigidly attached to the robotic end-effector, it exhibits a broader class of motions which can be used for parameter estimation [7].

However, using non-prehensile pushing to infer object parameters is a challenging task as the interaction dynamics between the object and the robotic system are sophisticated to model, due to uncertainty in the contacts, surface irregularities etc. [8], [9]. This can make it difficult to infer the parameters accurately. Therefore the data-driven approach is a more viable approach for such an interaction model.

To improve data efficiency and inference time, it is essential, the robotic systems, strategically select the next best exploratory push actions (active push exploration) such as where to push? and how to push? Authors have already shown that active object

A. Dutta and M. Kaboli are with the BMW Group, Munich, Germany. e-mail: name.surname@bmwgroup.com

A. Dutta and E. Burdet are with Imperial College of Science, Technology and Medicine, London, UK. M. Kaboli is with Radboud University, Netherlands.

Funded in part by the EU H2020 INTUITIVE under Grant ID 861166 and in part by EU Horizon PHASTRAC under Grant ID 101092096.

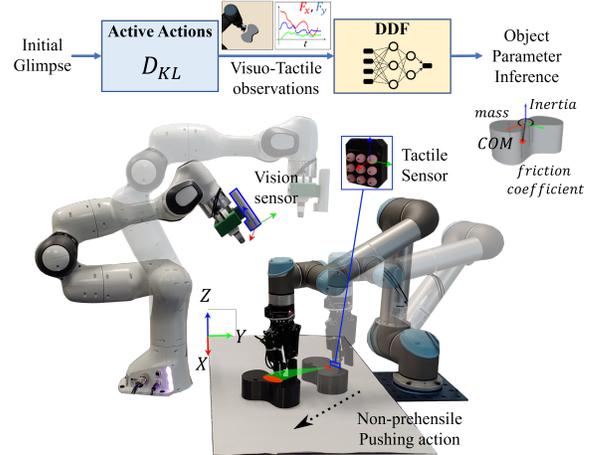


Fig. 1: Problem setup for visuo-tactile based active object parameter inference

exploration outperforms uniform and random strategy for reducing the uncertainty about objects for problems like object recognition [10]–[13] and pose estimation [14], [15]. In this work, we propose N -step Information Gain formulation for active exploratory push action selection which is crucial for dynamic interactions such as pushing.

Furthermore, since physical properties of objects are not directly observable and must be inferred from noisy visuo-tactile sensory observations, we introduce a novel dual differentiable filtering approach to address this and effectively handle the time-invariant nature of such properties.

II. RELATED WORKS

Estimating the physical properties of novel objects is a challenging problem in robotics, using either vision or tactile sensing. The physical object properties are not salient under static or quasi-static interactions, and often each parameter is only revealed under specific interactions, making it an interesting research problem [16].

One of the earliest works of Atkeson et al. [17] estimated the mass and moment of inertia of an object rigidly attached to a manipulator, using joint torques and a wrist-mounted force-torque sensor. Similar results have also been presented in [18]. These approaches required the object to be manually attached to the end-effector. Few works elevated this constraint as in [19], where authors used a 2-fingered bespoke mechanism to measure contact forces during planar pushing and in [20], the authors applied a tilting approach to measure wrenches to estimate the inertial parameters. Zhao et al. [21] incorporated friction estimation, by grasping the object and measuring the contact forces during the sliding regime. Most of these prior estimation techniques relied on precise force or tactile sensing, assumptions about the object geometry or the interaction between the object and environment, and employed specialized mechanisms, thus making it difficult for generalization and autonomous exploration of the object.

Some researchers attempted to overcome the previously mentioned limitations by introducing interactive manipulation techniques like grasping or pushing. In [22], the authors estimated only the mass of an object by controlled pushing, which required prior knowledge about the friction-co-efficient of the surface. Similarly, to determine the centre of mass of the object, Yao et al. [23] utilised the tactile forces during a 3-fingered robotic grasp. To estimate more object physical properties, Sundaralingam et al. [24] used a factor graph approach using in-hand manipulation with precise tactile and force-torque sensing. The approach relied on the approximation of in-hand object dynamics, known object shape and a marker-less tracking system. More recently, Uttayas et al. [25] estimated visco-elastic properties using filtering approach, relying on approximate spring mass damper model. The above-mentioned works employing interactive manipulation often used an analytical formulation to model the object-robot interaction, which is often approximate and has significant assumptions about the interactions.

Recently, data-driven and physics engine approaches are being taken to overcome such problems. Wu et al. [26] used deep learning to learn interactions between objects colliding in a physics engine and utilised the learned model to estimate mass and friction parameters for real object motion. Song et al. [27], [28] relied on a physics engine to predict expected object motions during pushing and employed Bayesian optimization on a real object motion to predict distributed mass and friction on objects. These works often relied on the accuracy of the physics engine and are generally computationally complex. Xu et al. [16] used only vision and deep learning to learn a representation of the mass and friction coefficient by randomly pushing and poking objects. Mavrakis et al. [29] collected large pushing trajectories (40k) in the simulation environment and learned a regression model for estimating an object’s inertial parameters during non-prehensile pushing. However, these approaches require intensive training and do not involve strategic interaction. In this work, we propose an active formulation for efficient training data-driven object-robot interaction model.

As of yet, either vision or tactile were used to estimate the physical object properties. On one hand, tactile information is crucial to infer multiple object properties like in [23], [24], [30], [31], however, requires precise position information and prior knowledge. On the other hand, vision-based approaches such as [16], [32] could only estimate fewer object properties with higher error rates, but required no prior knowledge about the object. To exploit the complementing vision and tactile sensing modalities, we propose to utilise both. Recently, Murali et al. [15] and Lee et al. [33] have shown visuo-tactile based approach significantly improves the performance of the robotics systems problems like pose estimation and contact-rich manipulation.

To tackle the above-mentioned problems and constraints for estimating or inferring the physical object properties, we propose visuo-tactile based active object parameter inference with a dual differentiable filter. Our proposed approach enables the robotic system to utilise vision and tactile information to interactively explore a novel object via non-prehensile pushing, not requiring any prior knowledge about the object’s shape, position etc. In addition, the novel N -step active formulation within the differentiable filtering facilitate efficient learning of the object-robot interaction model and also during inference by selecting the next best exploratory push actions.

Contributions

The contributions of the work are summarized in three folds.

- 1) We propose a novel dual differentiable filter to systematically account for the time-invariant nature of object parameters and time-varying object pose during non-prehensile pushing.
- 2) We propose a novel data-driven model for the dual differentiable filter which handles raw visuo-tactile sensory information and captures the object-robot interaction during pushing.
- 3) We propose a novel N -step lookahead formulation exploiting the differentiable filtering step to select the next best exploratory push actions for both learning the object-robot interaction model and performing inference.

We perform extensive simulation and real-robot experiments supporting the proposed method and comparison with the state-of-the-art baselines.

III. PROPOSED METHOD

A. Problem Definition

In this work, we consider the problem of estimating the state s of an unknown rigid object from vision o^V and tactile observation o^T using non-prehensile pushing actions a . At any given time t , the state $s_t = \{\psi_t, \phi\}$ comprises of time-varying factors: **pose**, $\psi_t = \{x_t, y_t, \theta_t\}$ as well as time-invariant factors: **parameters**, $\phi = \{m, \mu, CoM_x, CoM_y, I_z\}$ as *mass, relative friction coefficient between object and the surface, center of mass, inertia*. The center of mass is measured w.r.t frame attached to the geometric center of an object and only the rotational inertia in 2D is considered, as the interaction is restricted to motion in 2D. The observation o_t^V consists of RGB-D images of the pushing area and tactile observation o_t^T comprising of *2D contact forces, contact indicator*. The contact indicator $\in \{0, 1\}$, depending on whether the robot and object are in contact. The pushing action a_t is parameterized by *contact point (cp), push direction (pd)* and *velocity (v)* of the push. The *cp* comprises of 2D world co-ordinate of the contact point, *pd*, z-axis rotational angle of the robotic system aligned along a pushing direction & v is the magnitude of push velocity by the robotic system.

We perform quasi-static pushing [8] to infer the object parameters ϕ which are not directly observable either through vision or tactile sensing. Our proposed framework is illustrated in Fig.2 a). It comprises a novel dual differentiable filter for parameter and pose estimation along with data-driven models. The action selection for push affordance is done via computing N -step information gain term, making it an active dual differentiable filter (ADDF). Firstly, the robotic system learns data-driven models utilised within the differentiable filter. After learning, we perform inference on novel and unknown objects to estimate their parameters, with no prior information about the novel object. In the following sections, we explain the various components of the framework.

B. Bayesian inference and Differential Filters

We represent the belief about the current state of the object s_t with a distribution conditioned on previous actions $a_{1:t}$ and observations $o_{1:t}$. This distribution is denoted as the belief of the state $bel(s_t)$

$$bel(s_t) = p(s_t | o_{1:t}, a_{1:t}) = \frac{p(o_t | s_t, o_{1:t-1}, a_{1:t}) p(s_t | o_{1:t-1}, a_{1:t})}{p(o_t | o_{1:t-1}, a_{1:t})} \quad (1)$$

One prominent approach to computing the belief tractably is to employ Recursive Bayesian Filters which follow the structure as:

$$bel(s_t) = p(s_t | o_{1:t}, a_{1:t}) = \eta p(o_t | s_t, a_t) \overline{bel}(s_t) \quad (2)$$

$$\overline{bel}(s_t) = \int p(s_t | s_{t-1}, a_{t-1}) bel(s_{t-1}) ds_{t-1} \quad (3)$$

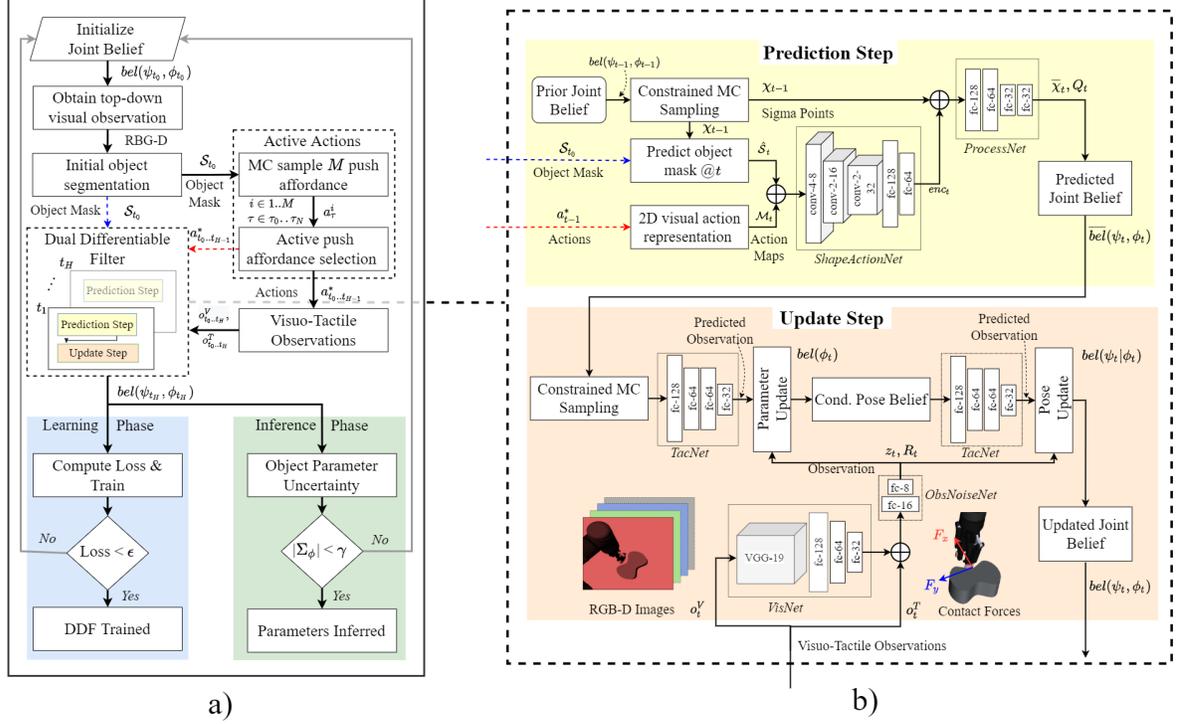


Fig. 2: Our proposed framework (ADDF) for visuo-tactile based active object exploration using non-prehensile manipulation. Part a) presents the overall framework and part b) presents an expanded view of the dual differentiable filter block.

Kalman Filters are a common choice of Bayesian Filtering which is optimal in linear systems and can be extended in non-linear cases through various approaches. Two key aspects of Bayesian filtering are the representation of the process model of the state in the form of $p(s_t|s_{t-1}, a_{t-1})$ and an observation likelihood model relating the states to the observations $p(o_t|s_t)$. For our problem, we employ a data-driven approach to learn the process and observation model along with the respective noise models, end-to-end using a differentiable filter. Recently, differentiable filters integrating Bayesian filtering with deep learning [34]–[37] were proposed. The authors have also shown that such an approach performs better compared to the standard deep-learning approach in handling real-world noise and in [38] showed the strength of such approach for a variety of tasks like visual optometry, visual object tracking etc.

In our problem, the data-driven models within a differentiable filter capture the complex and stochastic object-robot interaction model during non-prehensile pushing. Further, the pose of the object is intricately dependent on the parameters and straightforward combined (joint) filtering for pose and parameter does not perform well. We empirically show this in Section IV. Therefore, we utilize a dual filter design, exploiting the dependency among the states for consistent filtering and inferring the parameters of the object.

C. Dual Differentiable Filter

We derive our dual filter based on differentiable UKF [38], [39]. For the dual filter formulation, we explicitly represent the state s_t by the joint distribution of ψ_t and ϕ_t , via Multivariate Gaussian distribution:

$$bel(\psi_t, \phi_t) \sim \mathcal{N}(\psi_t, \phi_t | \mu_t, \Sigma_t) \quad (4)$$

with statistics $\mu_t \in \mathbb{R}^8$ and $\Sigma_t \in \mathbb{R}^{8 \times 8}$ as

$$\mu_t = \begin{pmatrix} \mu_{\psi_t} \\ \mu_{\phi_t} \end{pmatrix}, \quad \Sigma_t = \begin{pmatrix} \Sigma_{\psi_t} & \Sigma_{\psi_t \phi_t} \\ \Sigma_{\phi_t \psi_t} & \Sigma_{\phi_t} \end{pmatrix}. \quad (5)$$

The dual filter as shown in Fig.2(b) follows, the structure of a Kalman Filter with a *prediction step* and an *update step*, with key novelty explained in this section.

1) *Prediction Step*: In prediction step, the next step joint belief is predicted given the prior belief and the actions. The object parameters are real physical quantities with some physical constraints, (for e.g. $m, \mu > 0$). However, simply constraining the sigma points χ^{UT} in UKF approach does not preserve the true variance of the Gaussian distribution [40]. Therefore, we perform constrained Monte-Carlo sigma point sampling to preserve the physical constraints and the variance of the Gaussian. We employ a differentiable sampling method [41] to sample C sigma points on the joint distribution $bel(\psi_{t-1}, \phi_{t-1})$ instead of using standard Unscented Transform points:

$$\chi_{t-1}^{[i]} = \mu_{t-1} + \varepsilon^{[i]} \sqrt{\Sigma_{t-1}} \quad (6)$$

where, $i \in 1..C$ and $\chi_{t-1} = [\chi_{\psi_{t-1}}, \chi_{\phi_{t-1}}] \in \mathbb{R}^{C \times 8}$ with an associated weight $w_t^{[i]} = 1/C$ and $\varepsilon^{[i]} \sim \mathcal{N}(0, 1)$. We set $C = 100$ for all our experiments. The sigma points are filtered based on whether they satisfy the physical constraints and passed through the data-driven models. However, the invalid sigma points are also retained and reintroduced to preserve the uncertainty of the distribution. This is visually illustrated and explained in Appendix - Fig. 9.

Shape-Action Encoder: During the pushing, it is important to take into account the local geometry of the object and the action. Few previous works [32], [42] have shown that such an approach

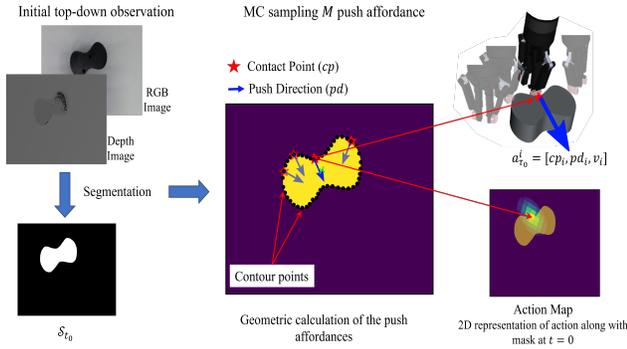


Fig. 3: Monte-Carlo (MC) action sampling and action maps illustration improves the action effect prediction. We encode the action along with the local geometry of the object at the point of contact to improve action effect predictions via the *ShapeActionEncoderNet*. *ShapeActionEncoderNet*: This comprises 3-layer CNN layers followed by 2 layers of feed-forward neural network. For each sigma point sampled from the current belief, an expected *segmentation mask* \hat{S}_t is generated by transforming the initial segmentation S_0 based on the pose information $\chi_{\psi_{t-1}}$. This represents the current geometry of the object at the point of action. Next, a 2D representation of the action - *action map* \mathcal{M}_t is generated. This is done via representing a 2D Gaussian distribution based on the action affordance $a_t = (cp, pd, v)$ in the image frame. Further details of the action map can be found in the appendix, and visually depicted in Fig.3. By this approach, we avoid generating complex object shape predictions for intricate visual perspective.

ProcessNet: The data-driven process model for predicting the change in *pose* of the object is approximated via 3 layer feed-forward neural network given prior joint sigma points and the shape-action encoding enc_t . In addition, we also employ the learnt heteroscedastic process noise model.

$$enc_{a_t} \leftarrow ShapeActionEncoderNet(\{\hat{S}_t, \mathcal{M}_t\}) \quad (7)$$

$$\bar{\chi}_{\psi_t}, Q_t \leftarrow ProcessNet(\chi_{t-1}, enc_{a_t}) \quad (8)$$

$$\bar{\chi}_{\phi_t} = \chi_{\phi_{t-1}} \quad (9)$$

where, $Q_t \in \mathbb{R}^{C \times 3}$ is the heteroscedastic diagonal covariance noise for time-varying pose. The predicted next step sigma points $\bar{\chi}_t$, along with the process noise Q_t are utilized to compute the expected Gaussian belief $\overline{bel}(\psi_t, \phi_t)$ as

$$\bar{\chi}_{\psi_t}^{[i]} = \bar{\chi}_{\psi_t}^{[i]} + \varepsilon^{[i]} \sqrt{Q_t^{[i]}} \quad (10)$$

$$\bar{\mu}_t = \sum_{i=0}^C w_t^{[i]} \bar{\chi}_t^{[i]} \quad (11)$$

$$\bar{\Sigma}_t = \sum_{i=0}^C w_t^{[i]} (\bar{\chi}_t^{[i]} - \bar{\mu}_t) (\bar{\chi}_t^{[i]} - \bar{\mu}_t)^T \quad (12)$$

where, $i \in 1..C$ and $\bar{\chi}_t = [\bar{\chi}_{\psi_t}, \bar{\chi}_{\phi_t}]$

2) *Update Step*: We recompute the constrained Monte-Carlo sigma point $\bar{\chi}_{\phi_t}$ sampling on the predicted belief $\overline{bel}(\psi_t, \phi_t)$ to incorporate the process noise. The dual filter employs a separate update of parameter belief similar to the parameter update presented in [43] and the conditional pose belief update based on the UKF update [39]. For updating the joint belief, we require an observation model to predict observation sigma points \bar{z}_t which has to account for both visual and tactile observations. To reduce the complexity of predicting raw RGB-D images, we split into the observation model two components, tactile and visual models. A *VisNet* network acts as a synthetic sensor generating the current noisy 2D pose

information x, y, θ from the current RGB-D images at each time. The *VisNet* comprises of first 10 layers of VGG-19 [44] pretrained on ImageNet followed by 3 layers of feed-forward network. For the tactile counterpart, a 4 layers of feed-forward network *TacNet* is utilised to predict the contact force information. In addition, a two-layer network *ObsNoiseNet* is also utilised to generate the heteroscedastic and diagonal observation noise.

$$\bar{z}_t^V = \bar{\chi}_{\psi_t}^V \bar{z}_t^T \leftarrow TacNet(\bar{\chi}_t^V, enc_{a_t}) \quad (13)$$

$$z_t^V \leftarrow VisNet(o_t^V), z_t^T = o_t^T \quad (14)$$

$$R_t \leftarrow ObsNoiseNet(z_t^V, z_t^T) \quad (15)$$

Parameter Update: We update the weights based on the likelihood of the observation sigma points $\bar{z}_t = [\bar{z}_t^T, \bar{z}_t^V]$ in the observation distribution $\sim \mathcal{N}(\cdot | z_t, Q_t)$

$$w_t^{[j]} = w_t^{[j]} e^{(-\frac{1}{2}(\bar{z}_t^{[j]} - z_t) R^{-1} (\bar{z}_t^{[j]} - z_t)^T)} \quad (16)$$

where $j \in 1..C$. The updated parameter belief $bel(\phi_t)$ is recomputed via a Gaussian Smooth Kernel [43] method after normalizing the updated weights.

$$\mu_{\phi_t} = \sum_{i=0}^C w_t^{[i]} \bar{\chi}_{\phi_t}^{[i]}; \quad m_{\phi_t}^{[i]} = a \bar{\chi}_{\phi_t}^{[i]} + (1-a) \mu_{\phi_t} \quad (17)$$

$$\Sigma_{\phi_t} = h^2 \sum_{i=0}^C w_t^{[i]} m_{\phi_t}^{[i]} - \mu_{\phi_t} \quad (18)$$

where a and $h = \sqrt{1-a^2}$ are shrinkage values of the kernels that are user-defined and set to 0.01, and m are the kernel locations.

Pose Update: We make use of the dependence of the pose on the parameters to compute the conditional pose distribution $bel(\psi_t | \phi_t) \sim \mathcal{N}(\psi_t | \mu_{\psi_t | \phi_t}, \Sigma_{\psi_t | \phi_t})$ using Multivariate Gaussian Theorem [45].

$$\mu_{\psi_t | \phi_t} = \psi_t + \Sigma_{\psi_t, \phi_t} \Sigma_{\phi_t}^{-1} (\phi_t - \mu_{\phi_t}) \quad (19)$$

$$\Sigma_{\psi_t | \phi_t} = \Sigma_{\psi_t} - \Sigma_{\psi_t, \phi_t} \Sigma_{\phi_t}^{-1} \Sigma_{\phi_t, \psi_t} \quad (20)$$

For conditional pose update, standard Unscented Kalman Filter (UKF) is employed on the predicted conditional pose distribution $\overline{bel}(\psi_t | \phi_t = \mu_{\phi_t})$ using Eq.20. The μ_{ϕ_t} of the updated parameter belief is utilised with predicted pose sigma points $\bar{\chi}_{\psi_t}^{UT}$ to obtain the predicted observation sigma points \bar{z}_t^V . The UKF update equations are skipped for brevity. After the conditional pose update, the posterior joint is computed as:

$$bel(\psi_t, \phi_t) = bel(\psi_t | \phi_t) bel(\phi_t) \quad (21)$$

Note, the cross-covariance matrices $\Sigma_{\psi_t, \phi_t}, \Sigma_{\phi_t, \psi_t}$ are not updated through the dual update step and are kept constant.

D. Active Non-prehensile Pushing Actions

The push action affordance is given by the tuple $a_t = (cp, pd, v)$. Possible *contact point's* cp and normal angle cn at the contact point is geometrically computed from the initial 2D segmentation S_0 based on our previous work in [15] and illustrated in Fig.3.

Monte-Carlo Sampling of push affordance: We generate M push affordances, $a_t^{[i]}$, $i \in 1..M$, from the possible points of contact points and contact normal by sampling a contact point and generating the $pd^{[i]} = cn^{[i]} + \delta$; $\delta \sim R(-5, 5)$ (deg). The velocity v is fixed for all cases keeping in mind the quasi-static assumption.

N-step Information Gain: To make the framework more sample efficient for real robot scenarios, we employ active action selection

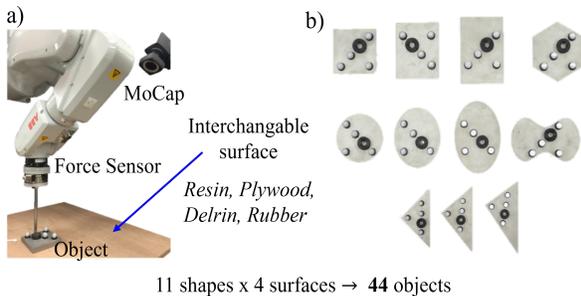


Fig. 4: MIT Push Dataset Setup [48]. Part a) presents the data collection setup. Part b) presents the various objects in the dataset

by formulating an N -step information gain criteria [46] under the filtering setting. We recursively use the prediction step of the dual differentiable filter without the update step to compute the expected Information Gain for both model learning and object parameter inference for each sampled non-prehensile pushing action $\pi^{[i]} = a_{\tau_0}^i \dots \tau_N$ over N -step in future $\tau = \tau_0 \dots \tau_N$

$$IG_N(\pi^{[i]}) \approx -\mathbb{E}_{p(\psi_{\tau_N}, \phi_{\tau_N} | \pi^{[i]})} [\ln(\overline{bel}^{[i]}(\psi_{\tau_N}, \phi_{\tau_N})) - \ln(\overline{bel}^{[i]}(\psi_{\tau_0}, \phi_{\tau_0}))] \quad (22)$$

where, $\overline{bel}^{[i]}(\psi_{\tau_N}, \phi_{\tau_N})$ is the hypothetical predictive joint distribution after N -step by taking action $\pi^{[i]}$ without taking account the actual observation. For our case, the expectation is computed as KL-Divergence form for which the closed form solution exists for Multivariate Gaussian distributions [47].

$$IG_N(\pi^{[i]}) \approx D_{KL}[\mathcal{N}^{[i]}(\psi_{\tau_N}, \phi_{\tau_N} | \overline{\mu}_{\tau_N}, \overline{\Sigma}_{\tau_N}) || \mathcal{N}^{[i]}(\psi_{\tau_0}, \phi_{\tau_0} | \overline{\mu}_{\tau_0}, \overline{\Sigma}_{\tau_0})] \\ \pi^* = \arg \max_{\pi^i} IG_N(\pi^{[i]}) \quad (23)$$

IV. EXPERIMENTS

In this section, we explain the experiment setup and the results obtained from the proposed method hereby referred to as VT-ADDF (Visuo-Tactile Active Dual Differentiable Filter). The closest state-of-the-art work to ours which dealt with the estimation of object parameters using robotic pushing was that of [29] and have taken it as the baseline. The baseline work utilised feature extraction using object pose, actions, and contact force information and a Multi-Output Regression Random Forest for data-driven regression modelling. We re-implemented the baseline approach to the best of our capability and validated the published results on the MIT Push Dataset.

In addition, we performed extensive ablation studies 1) exploring the efficiency of the active approach vs random and uniform actions for learning and inference, 2) employing only vision for parameter estimation under the dual filtering setup (termed at V-DDF Visual-Dual Differentiable Filtering). For this, the *TacNet* was removed and the observations were reduced to only RGB-D. The rest of the framework and dual filtering setup with the active actions remained the same. 3) Study of dual filtering approach compared to joint filtering (termed as VT-JDF). In this, instead of performing separate parameter and pose updates, only a single UKF update equation was used [38].

A. Experimental Setups

We tested our approach and compared the baseline on 3 experimental setups.

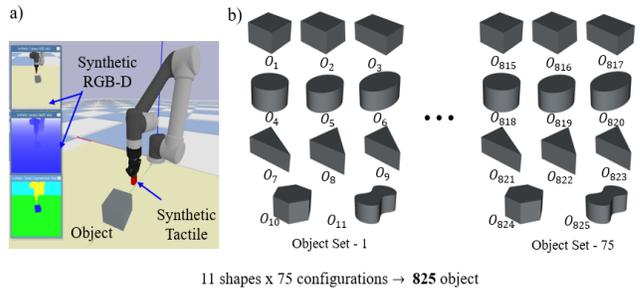


Fig. 5: Simulation Setup-*Sim Robotac*. Part a) presents the PyBullet scene of the setup. Part b) presents the object set (825) objects used for the experiments.

TABLE I: Parameter range for simulation setup

Property	Range of values
Mass (kg)	[0.2, 0.5, 0.8, 1.2]
μ	[0.35, 0.5, 0.7]
COM_x (m)	[-0.02, -0.015, 0, 0.015, 0.02]
COM_y (m)	[-0.02, -0.01, 0, 0.02, 0.03]
I_z (g.m ²)	[0.5, 0.9, 1.15, 1.5]
Shapes	11 (Fig. 5)

1) *Dataset - MIT Push Dataset*: We utilised the MIT Push dataset, a state-of-the-art robotic pushing dataset [49] with 44 different objects as shown in Fig.4. The dataset contains tactile, synthetic RGB-D, pose as well as parameter information. The objects were of 11 different shapes with varying mass, inertia, and 4 different surfaces - Abs, Delrin, Plywood, and Rubber-sheet were present. The center of mass of each object was slightly varied w.r.t its geometric center. The dataset has almost 10,000 pushes for each object, however, we selected a partial subset of pushes with no acceleration and velocity of 30 mm/s in a total of 3750 pushes. As this was a pre-recorded dataset, we only present estimation with uniform actions rather than active actions. We selected this setup to validate the baseline results as well as showcase that our proposed visuo-tactile dual differentiable filter can be utilized for different robotic environments.

2) *Simulation Setup - Sim Robotac*: We designed a simulation setup in PyBullet [50] to evaluate extensively our proposed approach as shown in Fig.5. In addition, it is possible to have a large set of objects with variations in physical parameters in the simulation which is often difficult real robotic setup. The setup comprised a simulated Robotiq gripper mounted on UR5 with a simulated tactile sensor attached to one of the finger pads. A synthetic RGB-D sensor was placed on top of the pushing area to simulate a visual sensor. In the simulation setup, 825 different objects were designed based on the parameters as presented in Table I. We utilised the simulation setup to perform extensive ablation studies, the results of which are presented in the following section.

3) *Robotic Setup*: The robotic setup consists of Universal Robots (UR5) augmented with Robotiq two-finger Gripper and a Panda robotic manipulator as shown in Figure 6. Tactile sensor [51] is attached to the outer surface of the finger of the gripper pads of the Robotiq Gripper and an Azure DK RGB-D camera is rigidly attached to the Panda Gripper. The maximum allowed speed for the UR5 was 25 mm/s for safety constraints. The ground truth values of the pose were collected using the motion capture system - Optitrack [52], whereas the ground truth values of the object parameters were computed from a CAD model of the objects. To obtain real objects with varying parameters, we designed configurable objects by 3D printing 4 shapes and adding additional weights at a precise

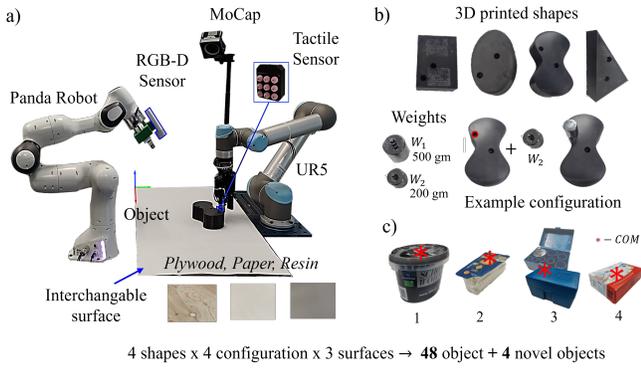


Fig. 6: Robotic Setup - Real Robotac. Part a) presents the robotic setup. Part b) presents an overview of the configurable objects. Part c) presents the novel objects selected to test

location in the objects, changing their mass, center of mass and inertia value. In addition, we utilised 3 different frictional surfaces - plywood, paper and resin sheet, to vary the relative friction coefficient between the object and the pushing surfaces. In total, we had 48 different objects after all possible configuration. In addition, we have 4 novel daily objects as shown in Figure.6 (c) which were not used in training and were kept only for testing. The objects had contrasting parameters (high mass of paint box (Object 1), high friction of sugar cube box (Object 4) and the plywood surface, as well as shifted COM in cheese (Object 2) and the weight box (Object 3).

B. Experimental Results

1) *Learning - DDF Training*: For training the networks used in the differentiable filter, we utilised a weighted combination of negative log-likelihood loss \mathcal{L}_{NLL} of the ground pose and parameter w.r.t to the belief and mean squared error loss \mathcal{L}_{MSE} of contact forces and synthetic pose. Iterative training was done using Adam optimizer till the loss converged.

For MIT Push Dataset, the time horizon was $t_H = 10s$ with a sampling rate of 10Hz. We split the 3750 trajectories, into 90% for training and 10% for inference. For the Sim-Robotac setup, the time horizon was $t_H = 15s$ with a sampling rate of 10Hz. 90% of 825 objects were utilised for training and 10% for testing, cross-validated 5 times. We performed an ablation study utilising a uniform, random and active approach of taking the action from the set of M -push affordances for training the filter. In addition, we also explored how much N -step lookahead is suitable. We chose N as 20%(=3 secs), 50% (=7.5 secs), and 70 (=10.5)% of the time horizon as future look-ahead steps for ablation study on active actions. We present the results of the ablation study in Fig. 7 (b), for learning efficiency. In addition, we also present the validation loss plots to highlight the stability and learning performance of the proposed approach (VT-ADDF) compared to using only vision (V-DDF) and utilizing a joint differentiable filter (VT-JDF) in Figure. 7 (a). For the Real-Robotac setup, the time horizon was $t_H = 10s$, with a sampling rate of 5 Hz and an active approach with 50% N -step lookahead (=5secs) selected for training the dual differentiable filters. 90% of the 48 3D printed configurable objects were utilised for training and 10% for testing.

C. Parameter Inference

For parameter inference of unknown (test) objects, we executed multiple push actions. At the end of each, the posterior belief of object parameters was utilised to initialize the belief for the next push. We present the results of the parameters $m, \mu, CoM_x, CoM_y, I_z$

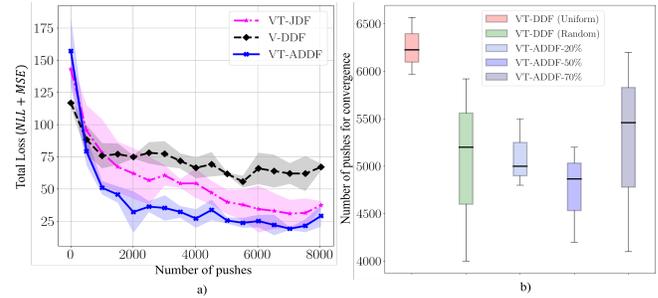


Fig. 7: Learning results of the ablation studies in Sim-Robotac setup. Part (a) presents comparative performance learning stability of VT-ADDF vs V-DDF vs VT-JDF. Part (b) presents the learning efficiency of different push action selection methods - Uniform, Active, Random

inference for the ablation study in Sim-Robotac setup in Table.II. For parameter inference, the N -step lookahead was selected as 50% of time horizon t_H for both Sim-Robotac and Real-Robotac setups. In addition, Fig.8 presents a closer look at the filtering action of the different ablation approaches and parameter inference convergence in the Sim-RoboTac setup. Further, the comparative parameter estimation results of the proposed approach (VT-ADDF) compared to the baseline work of [29] are presented in Table.III. We present separate estimation results for the novel objects which were not utilised for training, to evaluate the generalisation of the proposed approach compared to the baseline. As the range of values of different parameters like mass, friction co-efficient and inertia are different and the values often close to 0 (for the center of mass) we employ a normalised root mean square $NRMSE$ [29] as metric to evaluate the performance over different parameters which is the root mean squared error divided by the range of values ($\psi_{max} - \psi_{min}$) of each parameters in the setups. Lower value, signifies better estimation.

D. Discussion

In this work, we proposed a novel visuo-tactile based active object parameter inference with a dual differentiable filter.

The results of the ablation study for learning show that active actions significantly improve the sample efficiency by around 20% (push actions) compared to uniform actions and has lower variance than random action selection as presented in Fig.7(b). Moreover, it is shown that the dual filtering approach has more stable learning than the joint filter as presented in Fig.7(a). This demonstrates the efficacy of our proposed novel dual differentiable filtering approach compared to a joint differentiable filter method. Furthermore, our experimental results show that by using only vision, the network fails to reduce the loss and tends to overfit. This is expected, as the parameter estimation is difficult only via vision and has high error rates, leading to higher loss values.

The obtained results from the parameter inference show that our proposed approach performs consistently on different experimental setups - MIT Push Dataset, Sim-Robotac and Real-Robotac, compared to baseline work of [29], which fails to generalise for novel objects in real robotic setup. Moreover, the limitation of providing ground truth pose information in the baseline approach is elevated by our proposed ADDF approach during the inference step. Furthermore, the ablation study shows that active actions have a better estimation of parameters compared to uniform and with lower variance than random actions with the same number of push actions. Compared to using only vision, the visuo-tactile dual differentiable approach performs much better, especially in parameters like the center of mass prediction, as well as, is more stable and accurate than the joint filtering approach. Through the different setups, we

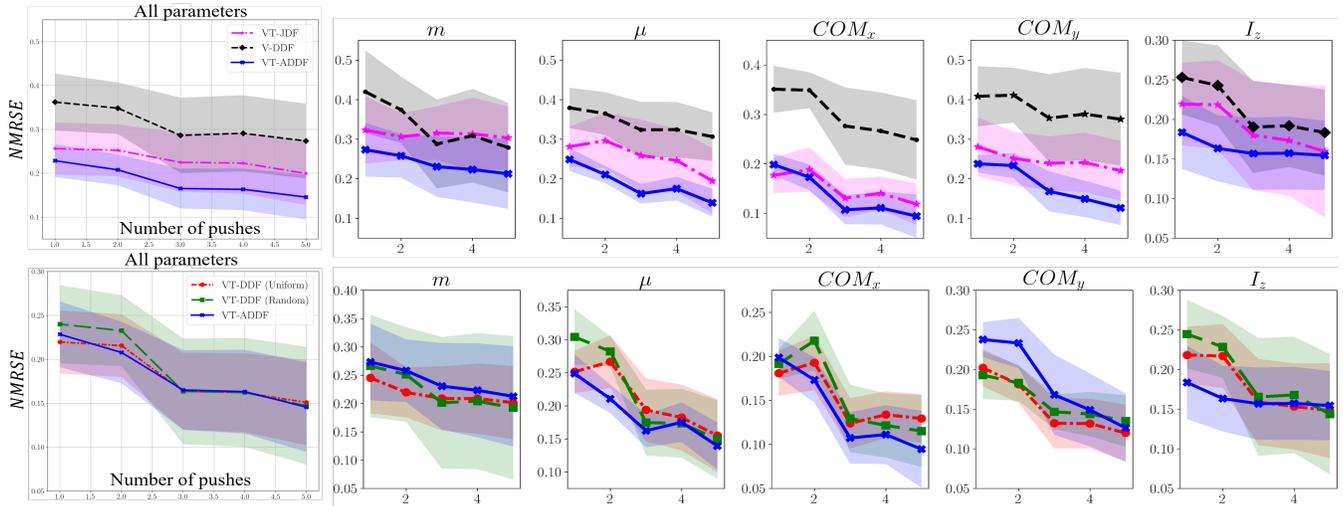


Fig. 8: Inference result during the filtering step presented after each push action.

TABLE II: Inference result of $NRMSE$ values of different parameters in the ablation study in Sim-Robotac setup

	$mass$	μ	com_x	com_y	I_z	Overall
VT-JDF	0.33 ± 0.08	0.19 ± 0.08	0.12 ± 0.04	0.22 ± 0.08	0.16 ± 0.08	0.20 ± 0.07
V-DDF	0.28 ± 0.11	0.31 ± 0.06	0.25 ± 0.08	0.35 ± 0.12	0.18 ± 0.05	0.27 ± 0.09
VT-ADDF	0.21 ± 0.09	0.14 ± 0.04	0.09 ± 0.04	0.13 ± 0.04	0.15 ± 0.03	0.14 ± 0.05
VT-DDF (Uniform)	0.20 ± 0.06	0.16 ± 0.06	0.13 ± 0.03	0.12 ± 0.03	0.12 ± 0.06	0.15 ± 0.05
VT-DDF (Random)	0.19 ± 0.13	0.15 ± 0.06	0.11 ± 0.04	0.13 ± 0.03	0.14 ± 0.08	0.14 ± 0.07

TABLE III: Parameter Inference result of $NRMSE$ value for proposed approach VT-ADDF compared to baseline work of [29] in various setups

<i>Experimental Setup</i>	$mass$		mu		com_x		com_y		I_z	
	Baseline	VT-ADDF								
MIT Push Dataset	0.11 ± 0.1	0.19 ± 0.02	0.18 ± 0.04	0.17 ± 0.02	0.13 ± 0.06	0.10 ± 0.04	0.12 ± 0.09	0.09 ± 0.07	0.17 ± 0.02	0.16 ± 0.01
Sim Robotac	0.14 ± 0.06	0.21 ± 0.09	0.16 ± 0.06	0.14 ± 0.04	0.18 ± 0.12	0.09 ± 0.04	0.14 ± 0.15	0.13 ± 0.04	0.20 ± 0.11	0.15 ± 0.03
Real Robotac (RR)	0.25 ± 0.12	0.22 ± 0.09	0.14 ± 0.03	0.19 ± 0.06	0.12 ± 0.08	0.10 ± 0.01	0.20 ± 0.1	0.11 ± 0.07	0.16 ± 0.05	0.09 ± 0.01
RR Novel Objects	0.29 ± 0.09	0.20 ± 0.05	0.21 ± 0.03	0.18 ± 0.06	0.17 ± 0.09	0.11 ± 0.05	0.22 ± 0.05	0.12 ± 0.05	0.22 ± 0.10	0.15 ± 0.08

also show that the proposed visual tactile-based dual differentiable filter for parameter inference is agnostic to robotic setups as long as sufficient visual and tactile information is present.

From the limitations, our proposed ADDF requires separate training of the *VisNet* with *MSE* loss to obtain pose information for the novel objects which are visually quite different from the training set. Instead of using RGB-D as visual observations and a 2D pose estimation network, it will be viable to use point clouds and avoid the requirement of using pose estimation altogether or use recent one-shot pose-estimation approaches. In addition, it will be interesting to avoid the requirement of having ground truth states and parameter values during training as well as develop a framework which can discover physical object representations.

V. CONCLUSION

In this work, we addressed the problem of estimating the properties of rigid objects using vision and tactile observations solely via non-prehensile pushing. The proposed approach first learns an object interaction model using known objects, which is utilized for inference of novel objects under differentiable filter settings. We present a novel formulation of active action selection with the differentiable filter as one of the key contributions. The generalizable capability of the framework makes it viable for real robotic applications and opens the possibility to explore the approach for other interaction techniques for object parameter estimation like grasping.

ACKNOWLEDGMENT

We would like to thank Dr. Alina Kloss for her useful insights and open sourced implementation of the differentiable filters. In addition, we would like to thank Prajval Kumar Murali, Iman Nematollahi, and Dr. Xiaoxiao Cheng for their constructive reviews.

REFERENCES

- [1] E. Burdet *et al.*, *Human Robotics: neuromechanics and motor control*. MIT press, 2013.
- [2] M. Kaboli *et al.*, “Tactile-based manipulation of deformable objects with dynamic center of mass,” in *JCHR*. IEEE, 2016.
- [3] J. Bohg *et al.*, “Interactive perception: Leveraging action in perception and perception in action,” *IEEE Trans. on Rob.*, vol. 33, no. 6, pp. 1273–1291, 2017.
- [4] R. Bajcsy *et al.*, “Revisiting active perception,” *AuRo*, vol. 42, 02 2018.
- [5] L. Seminara *et al.*, “Active haptic perception in robots: a review,” *Front. in Neurobotics*, vol. 13, p. 53, 2019.
- [6] Q. Li *et al.*, “A review of tactile information: Perception and action through touch,” *IEEE Trans. on Rob.*, vol. 36, no. 6, pp. 1619–1634, 2020.
- [7] M. T. Mason, “Progress in nonprehensile manipulation,” *IJRR*, vol. 18, no. 11, pp. 1129–1141, 1999.
- [8] M. T. Mason, “Mechanics and planning of manipulator pushing operations,” *IJRR*, vol. 5, no. 3, pp. 53–71, 1986.
- [9] J. Stüber *et al.*, “Let’s push things forward: A survey on robot pushing,” *Front. in Robotics and AI*, p. 8, 2020.
- [10] M. Kaboli *et al.*, “Tactile-based active object discrimination and target object search in an unknown workspace,” *AuRo*, vol. 43, pp. 123–152, 2019.

[11] M. Kaboli *et al.*, “Active tactile transfer learning for object discrimination in an unstructured environment using multimodal robotic skin,” *IJHR*, vol. 15, no. 01, p. 1850001, 2018.

[12] M. Kaboli *et al.*, “A tactile-based framework for active object learning and discrimination using multimodal robotic skin,” *IEEE RAL*, vol. 2, no. 4, pp. 2143–2150, 2017.

[13] D. Feng *et al.*, “Active prior tactile knowledge transfer for learning tactual properties of new objects,” *Sensors*, vol. 18, no. 2, p. 634, 2018.

[14] Murali, P. K. *et al.*, “Active visuo-tactile point cloud registration for accurate pose estimation of objects in an unknown workspace,” in *IEEE IROS*, 2021, pp. 2838–2844.

[15] Murali, P. K. *et al.*, “Active visuo-tactile interactive robotic perception for accurate object pose estimation in dense clutter,” *IEEE RAL*, vol. 7, no. 2, pp. 4686–4693, 2022.

[16] Z. Xu *et al.*, “Densephysnet: Learning dense physical object representations via multi-step dynamic interactions,” *arXiv preprint arXiv:1906.03853*, 2019.

[17] C. G. Atkeson *et al.*, “Estimation of inertial parameters of manipulator loads and links,” *IJRR*, vol. 5, no. 3, pp. 101–119, 1986.

[18] C. Wang *et al.*, “Parameter estimation and object gripping based on fingertip force/torque sensors,” *Measurement*, vol. 179, p. 109479, 2021.

[19] Y. Yu *et al.*, “Estimation of object inertia parameters on robot pushing operation,” in *ICRA*. IEEE, 2005, pp. 1657–1662.

[20] Y. Yu *et al.*, “Estimation of mass and center of mass of grasplless and shape-unknown object,” in *ICRA*, vol. 4. IEEE, 1999, pp. 2893–2898.

[21] Z. Zhao *et al.*, “Center of mass and friction coefficient exploration of unknown object for a robotic grasping manipulation,” in *IEEE ICMA*, 2018, pp. 2352–2357.

[22] S. Tanaka *et al.*, “Active mass estimation with haptic vision,” in *IEEE ICPR*, vol. 3. IEEE, 2004, pp. 256–261.

[23] K. Yao *et al.*, “Tactile-based object center of mass exploration and discrimination,” in *IEEE ICHR*, 2017, pp. 876–881.

[24] B. Sundaralingam and T. Hermans, “In-hand object-dynamics inference using tactile fingertips,” *IEEE Trans. on Robo.*, vol. 37, no. 4, pp. 1115–1126, 2021.

[25] P. Uttayopas, “Object recognition using mechanical impact, viscoelasticity, and surface friction during interaction,” *IEEE Trans. on Haptics*, vol. 16, no. 2, pp. 251–260, 2023.

[26] J. Wu *et al.*, “Galileo: Perceiving physical object properties by integrating a physics engine with deep learning,” *Advances in Neural Info. Processing Systems*, vol. 28, 2015.

[27] C. Song and A. Boularias, “A probabilistic model for planar sliding of objects with unknown material properties: Identification and robust planning,” in *IEEE IROS*. IEEE, 2020, pp. 5311–5318.

[28] C. Song and A. Boularias, “Learning to slide unknown objects with differentiable physics simulations,” in *RSS*, 2020.

[29] N. Mavrakis *et al.*, “Estimating an object’s inertial parameters by robotic pushing: a data-driven approach,” in *IEEE IROS*. IEEE, 2020, pp. 9537–9544.

[30] M. Kaboli *et al.*, “Humanoids learn touch modalities identification via multi-modal robotic skin and robust tactile descriptors,” *Advanced Robotics*, vol. 29, no. 21, pp. 1411–1425, 2015.

[31] M. Kaboli *et al.*, “In-hand object recognition via texture properties with robotic hands, artificial skin, and novel tactile descriptors,” in *ICHR*. IEEE, 2015, pp. 1155–1160.

[32] A. Kloss *et al.*, “Accurate vision-based manipulation through contact reasoning,” in *IEEE ICRA*. IEEE, 2020, pp. 6738–6744.

[33] M. A. Lee *et al.*, “Making sense of vision and touch: Self-supervised learning of multimodal representations for contact-rich tasks,” in *IEEE ICRA*. IEEE, 2019, pp. 8943–8950.

[34] R. Jonschkowski *et al.*, “Differentiable particle filters: End-to-end learning with algorithmic priors,” *arXiv preprint arXiv:1805.11122*, 2018.

[35] K. T. Yu and A. Rodriguez, “Realtime state estimation with tactile and visual sensing, application to planar manipulation,” in *IEEE ICRA*. IEEE, 2018, pp. 7778–7785.

[36] A. S. Lambert *et al.*, “Joint inference of kinematic and force trajectories with visuo-tactile sensing,” in *IEEE ICRA*. IEEE, 2019, pp. 3165–3171.

[37] T. Haaroja *et al.*, “Backprop kf: Learning discriminative deterministic state estimators,” *Advances in Neural Info. Process. Systems*, vol. 29, 2016.

[38] A. Kloss *et al.*, “How to train your differentiable filter,” *AuRo*, vol. 45, no. 4, pp. 561–578, 2021.

[39] S. Thrun *et al.*, “Probabilistic robotics,” *Communications of the ACM*, vol. 45, no. 3, pp. 52–57, 2002.

[40] D. Ebeigbe *et al.*, “A generalized unscented transformation for probability distributions,” *ArXiv*, 2021.

[41] M. Wüthrich *et al.*, “Robust gaussian filtering using a pseudo measurement,” in *2016 American Control Conference (ACC)*. IEEE, 2016, pp. 3606–3613.

[42] G. Ziyan *et al.*, “Planar pushing of unknown objects using a large-scale simulation dataset and few-shot learning,” in *IEEE CASE*. IEEE, 2021, pp. 341–347.

[43] J. Liu and M. West, “Combined parameter and state estimation in simulation-based filtering,” in *Sequential Monte Carlo methods in practice*. Springer, 2001, pp. 197–223.

[44] K. Simonyan and A. Zisserman, “Very deep convolutional networks for large-scale image recognition,” *arXiv preprint arXiv:1409.1556*, 2014.

[45] C. B. Do, “The multivariate gaussian distribution,” *Section Notes, Lecture on Machine Learning, CS*, vol. 229, 2008.

[46] K. Friston, “The free-energy principle: a unified brain theory?” *Nature Reviews neuroscience*, vol. 11, no. 2, pp. 127–138, 2010.

[47] J. Duchi, “Derivations for linear algebra and optimization,” *Berkeley, California*, vol. 3, no. 1, pp. 2325–5870, 2007.

[48] “The mcube lab - push dataset,” <https://mcube.mit.edu/push-dataset/index.html>, (Accessed on 03/02/2023).

[49] K. Yu *et al.*, “More than a million ways to be pushed. a high-fidelity experimental dataset of planar pushing,” in *IEEE IROS*. IEEE, 2016, pp. 30–37.

[50] E. Coumans and Y. Bai, “Pybullet, a python module for physics simulation for games, robotics and machine learning,” 2016.

[51] Contactile, “Contactile,” <https://contactile.com/>, 2022, [Online; accessed 15092022].

[52] “Optitrack - motion capture systems,” <https://optitrack.com/>, (Accessed on 03/02/2023).

APPENDIX

Constrained Monte Carlo Sampling

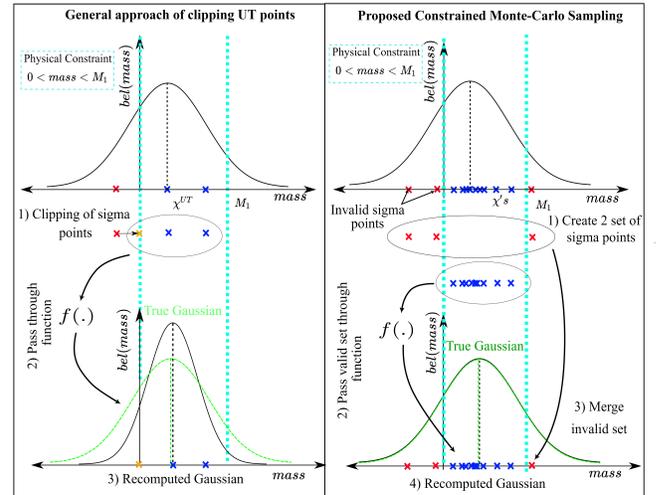


Fig. 9: Constrained Monte Carlo sampling

Action Maps A 2D Gaussian can be generated by the following equations, along the push direction pd and centered at the contact point in the image frame (cp_{if}).

$$px = \begin{pmatrix} pixel_x \\ pixel_y \end{pmatrix}, K = \begin{pmatrix} \frac{\cos^2(pd)}{2v^2} + \frac{\sin^2(pd)}{2} & \frac{\sin(2pd)}{4v^2} - \frac{\sin(2pd)}{4} \\ \frac{\sin(2pd)}{4v^2} - \frac{\sin(2pd)}{4} & \frac{\sin^2(pd)}{2v^2} + \frac{\cos^2(pd)}{2} \end{pmatrix}$$

$$\mathcal{M}_t = e^{(-\frac{1}{2}(px - cp_{if})K(px - cp_{if})^T)} \quad (24)$$