

Soft-bubble grippers for robust and perceptive manipulation

Naveen Kuppuswamy, Alex Alspach, Avinash Uttamchandani, Sam Creasey, Takuya Ikeda, and Russ Tedrake*

Abstract—Manipulation in cluttered environments like homes requires stable grasps, precise placement and robustness against external contact. We present the *Soft-bubble* gripper system with a highly compliant gripping surface and dense-geometry visuotactile sensing, capable of multiple kinds of tactile perception. We first present various mechanical design advances and a fabrication technique to deposit custom patterns to the internal surface of the sensor that enable tracking of shear-induced displacement of the manipulant. The depth maps output by the internal imaging sensor are used in an in-hand *proximity* pose estimation framework - the method better captures distances to corners or edges on the manipulant geometry. We also extend our previous work on tactile classification and integrate the system within a robust manipulation pipeline for cluttered home environments. The capabilities of the proposed system are demonstrated through robust execution multiple real-world manipulation tasks. A video of the system in action can be found [here](#).

Index Terms—soft-robotics, robot manipulation, tactile sensing, shear sensing, visuotactile

I. INTRODUCTION

Incorporating mechanical compliance, or softness, into manipulators has been observed to be beneficial in making grasping more robust [1]. However, many in-home tasks necessitate not just robust picking, but precise placement, e.g. dishwasher loading, liquid container handling, and shelf stacking. Although compliance can somewhat mitigate inaccuracies in pre-grasp pose estimation and object shape variations, variability in the post-grasp state may adversely impact task success. This problem is further exacerbated in highly cluttered home environments with visual occlusions, tight spatial constraints, and unforeseen contacts. Augmenting compliant end effectors with tactile perception capabilities might help compensate for variability in the post-grasp state of manipulands and for unexpected disturbances. To address these challenges, we present the highly compliant *Soft-bubble* gripper system which integrates multiple tactile perception capabilities in order to enable robust manipulation in tightly constrained environments.

While progress has been made in tactile sensing techniques, open questions remain on the ideal design for soft or compliant surfaces as well as on what kind of perception they should enable. For manipulation in cluttered environments, tactile sensing should enable at least a few of the following capabilities: manipulant pose estimation (to estimate post-grasp pose), manipulant external force estimation (to estimate both expected and unexpected forces), slip detection, object classification (e.g. to compensate for occlusion when grasping

in clutter), and grasp quality estimation. This paper tackles several of these problems.

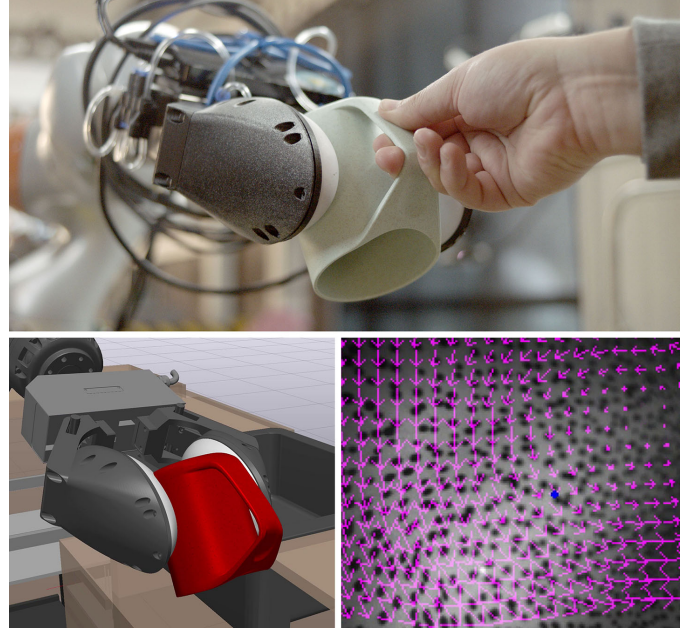


Fig. 1: Top: Highly compliant *Soft-bubble* parallel gripper on a robot arm grasping an object. Bottom-left: In-hand pose estimate during this interaction; Bottom-right: Shear-induced displacements tracked on the interior of this sensor.

The *Soft-bubbles* [2], [3] combine a highly compliant manipulation surface with dense geometry sensing. Mechanically, they achieve robust grasps since they are highly deformable, are easy to build due to their air-filled membrane design, and are fairly durable. They are also closely related to other visuotactile sensors reported in literature [4], with the key difference that the generated depth maps are directly measured by the internal imaging sensor and are not inferred. Some notable limitations on the prior iterations of the *Soft-bubbles* were their large form factor and an inability to track shear forces on the contact surface. This meant they could neither be used for human-scale object manipulation nor compensate for externally induced shear. Meaningfully co-designing the mechanical and perceptual aspects of these sensors for a robust manipulation pipeline was an open question - one which is directly tackled in this paper.

We present four key contributions to the *Soft-bubble* technology that combine to allow for robust manipulation of real-world, human-scale objects in a cluttered home environment:

(A) Design and fabrication techniques for a smaller parallel gripper form factor, leveraging modifications of the previously used time-of-flight (ToF) camera system

* All authors are with the Toyota Research Institute, One Kendall Square, Building 600, Cambridge, MA 02139, USA, [firstname.lastname]@tri.global

- (B) Techniques for adding high-density internal markers to the bubble surface that enable dense optical flow pattern tracking for estimating shear-induced membrane displacement
- (C) An optimization-based *proximity* pose estimation framework that is inspired by current state-of-the-art depth tracking methods.
- (D) Integrated tactile classification, using simultaneous images and an automated ground-truth labeling pipeline

We demonstrate the capabilities of the *Soft-bubble* gripper system through various tasks that test its sensing effectiveness, robustness to uncertainty, as well as its efficacy in human-robot interaction and antagonism. We conclude with observations on the performance, the gripper’s limitations, and a discussion of future work.

II. RELATED WORK

The contributions presented in this paper cover four key areas and the related work for each is summarized in this section.

A. Compliant Gripper Design

In contrast to most compliant gripper designs, which are based on solid materials [1], the *Soft-bubbles* use an air-filled deformable structure such as those reported in Kim et al [5]. They are also a high-resolution visuotactile sensor, like GelSight [6], [7], and GelSlim [8], which is advantageous for precision manipulation. The air-inflated latex membrane structure of the *Soft-bubble* has multiple mechanical benefits. (i) a large degree of compliance through a combination of pneumatic effects and membrane elasticity¹, (ii) the resulting higher friction at the contact patch leads to better grasps, (iii) they are potentially more robust to wear-and-tear from long-term use.

B. Shear Estimation

Estimating shear and/or slip from visuotactile sensing has been investigated [9], particularly for ensuring stable and sensitive grasps that are robust to unforeseen external contacts. A popular approach has been to embed uniform visual markers either on the contacting surface or within the the sub-surface material. Methods have also been proposed for decomposing this shear information into tangential and torsional components [10].

However, due to the large deformations and ellipsoidal form factor of the *Soft-bubbles*, such uniform markers result in shear tracking errors. Therefore we developed a pseudorandom pattern deposition technique that enables dense optical flow processing for shear-deformation tracking.

¹Typically, the deformations exhibited by *Soft-bubbles* are at least an order of magnitude greater than GelSight or GelSlim (multiple cm vs multiple mm on human-scale objects like mugs)

C. In-hand Pose Estimation

The problem of estimating the post-grasp pose of a manipulator within an end effector’s grasp has been attempted both with and without the use of tactile sensors. State-of-the-art methods for tracking known object geometries using depth sensors typically leverage probabilistic filters that account for unforeseen occlusions [11]. One application of tactile sensing has been to improve the stability of the filter under contact by collapsing the state to a contact manifold [12]. Algorithms such as DART overcome some of the intractability issues in filter-based tracking by solving an optimization problem and can also deal with articulated objects [13]; more accurate extensions have been proposed that account for physical non-penetration constraints when in contact with known surfaces [14].

In the case of visuotactile sensors, the rendered depth maps can be of sufficiently high resolution to be used directly for in-hand pose estimation [15]. Direct depth registration methods such as Iterative Closest Point (ICP) can be employed for local pose estimation [2] given a sufficiently close initial guess that can either be learned [16] or be obtained from methods that can compute the contact-patch resulting from contact with a tactile sensor [3]. Also, external depth sensor information can be fused with tactile depth [17].

The *proximity* pose estimator we present in this paper is closely related to optimization-based methods like DART using visuotactile sensors [17] due to their robustness against initialization - we propose some strategies to improve the smoothness of the cost function. The deformations of the *Soft-bubbles* under stable grasp capture enough of the object geometry to enable an optimization-based pose estimator from tactile information alone, as presented in this paper.

D. Tactile Classification

Tactile material and object classification is an interesting use case for tactile sensors [18] since it enables recognition of the object that is being grasped even when it is occluded from other vision sensing. As in [2], we use a ResNet18 architecture [19] image classifier; we now train and infer on concatenated images simultaneously combining individual depth or IR images from each of the sensors. An automatic labelling pipeline was developed for generating sufficient quantities of labelled images.

III. MECHANICAL DESIGN

Our previous *Soft-bubble* prototype was relatively large with a single sensor used as an end effector [2]. The new *Soft-bubble* fingers, as seen in Fig. 2, have been designed to both task-based and perceptual requirements. In order to achieve tasks in constrained domestic environments, the *Soft-bubbles* are designed to be attached to a standard parallel gripper, to interact with human-scale objects, and to fit into tight household spaces (e.g. a sink or dish washer). Perceptually driven improvements include the use of a shorter range depth sensor, interchangeable “bubble modules,” as well as a method

for adding visually trackable patterns to a bubble's internal surface.

A. Internal Depth Sensor

For shorter-range depth sensing (as compared to the off-the-shelf picoflexx), a prototype ToF depth sensor from PMD with a shorter imager-emitter baseline was used. This shorter baseline allows for more consistent illumination of the deformable bubble membrane at closer distances, and therefore permits more accurate depth measurements. As seen in Fig. 3, the internal ToF sensor is angled relative to the contact surface, with a working range of 4-11 cm. Depth measurements are more accurate near the upper bound of this range. The angled ToF depth sensor maximizes field of view (FOV) and depth measurement accuracy at the center and tips of the fingers where contact happens most, while minimizing the overall gripper width.

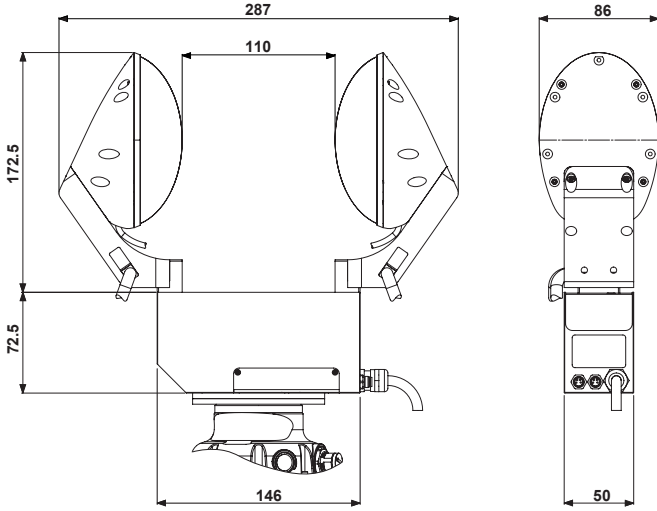


Fig. 2: The *Soft-bubble* gripper is designed to manipulate household objects in tight spaces. All dimensions in mm.

B. Bubble Modules

A modular design allows for quick swapping of bubbles. Bubbles may be exchanged for various reasons including the replacement of a damaged module, switching mechanical properties like membrane thickness, or swapping internal patterning based on the needs of the perception system. The modules consist of a plain or patterned latex membrane, a clear acrylic backing plate, a ring that mounts the membrane onto the acrylic, threaded inserts, and a tube fitting for inflation and pressure sensing. An exploded view of this cartridge can be seen in Fig. 4. The components are sealed together using CA glue. The sensor electronics remain mounted to the gripper, while the bubble cartridge can be removed and replaced using four screws.

C. Printed Patterning

The various patterns seen in Fig. 5 are script-generated with parameters controlling dot density (dots/mm²), minimum

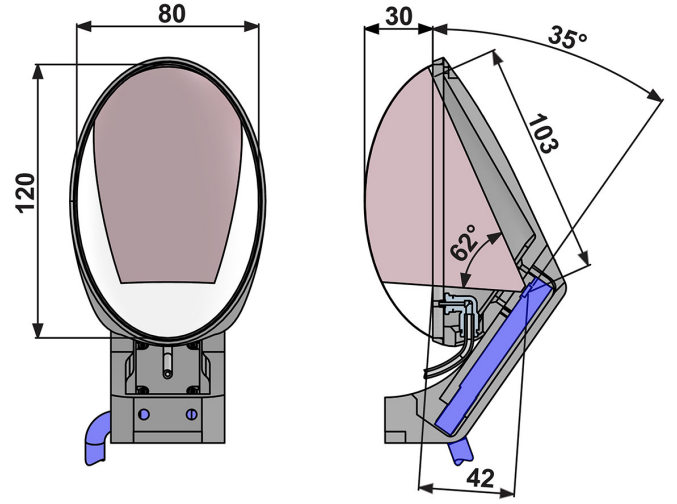


Fig. 3: An angled, short-range ToF depth sensor (blue) maximizes FOV (red) and depth measurement accuracy while minimizing overall gripper width. All dimensions in mm.



Fig. 4: Assembled and exploded views of bubble module components including latex membrane, clear acrylic plate, ring that mounts the membrane onto the acrylic, threaded inserts, tube fitting and tube. The red surface is where CA glue is applied during assembly to seal the module.

and maximum dot diameter, and pattern randomness. These patterns are output as vector files for laser cutting adhesive-backed stencils. These stencils are adhered to plain sheet latex, a thin layer of balloon screen printing ink is painted over the stencil, then the stencil is removed while the ink is still wet. This method can be used to apply patterns of single and multiple colors.

D. Construction and Durability

A batch of *Bubble-sensor* fingers can be inexpensively assembled in as little as two hours with no more than a FDM printer, laser cutter, scissors, glue and a paint brush; there is no casting, as we use off-the-shelf sheet latex for the membrane. In long-running experiments, we were able to operate the



Fig. 5: Experimental bubble modules and stencils with a variety of pseudorandom printed internal dot patterns for shear tracking, with dots in both black and gray.

Soft-bubbles for over 100 hours before the first failure - a puncture from the sharp edge of a manipulant. This type of failure is simple to recover from. When a puncture occurs, the surrounding area is lightly sanded then cleaned. A small patch is cut from latex, then both the patch and area surrounding the puncture are thinly painted with rubber cement, allowed to dry for a few seconds, then the patch applied. After a few minutes, the bubble can be re-inflated with negligible effect on the associated tactile sensing capabilities. In the event that a bubble can not be patched, the module is then swapped out for a backup bubble. We are currently exploring techniques to improve the durability and the mean time to failure.

IV. PROXIMITY POSE ESTIMATION

As mentioned in Sec. II, our approach to pose estimation of known objects is closely related to optimization-based methods for tracking on point-clouds [13], [14], [17] with a few modifications. In our prior work [2], [3], we observed that quality and resolution of the pair of images produced by the *Soft-bubbles* on contact are more than sufficient to enable tracking of manipulant pose. Here, we focus on a method that is more robust to estimator initialization - this is vital in the case of manipulating in cluttered environments where a priori guesses on in-hand state may be poor or non-existent.

Consider the scenario presented in Fig. 6. There is a rigid object O with an associated geometry reference frame G that is held between two *Soft-bubbles* mounted on the fingers of a robot gripper. There is a frame C_i associated with each of the *Soft-bubbles* and the gripper in turn has a tool frame T associated with it. The single-shot in-hand pose estimation problem is to estimate ${}^G X_T(t)$, for a known object geometry (where ${}^G X_T \in \text{SE}(3)$), when the object is under a stable grasp (the object is not slipping) using available sensor measurements $Y(t)$.

As previously suggested by Schmidt et al [13], signed distance functions (SDF) are a convenient parameterization of the object geometry since they can overcome the difficulty

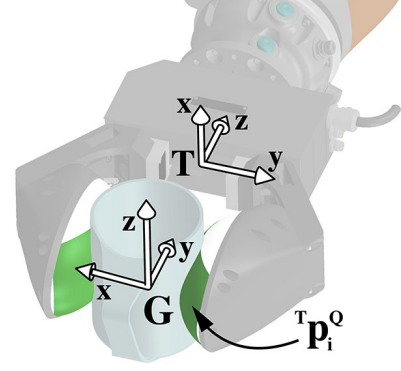


Fig. 6: The end-effector and manipulant geometry frames for the in-hand pose estimation problem. The shaded region in green represents the point-cloud corresponding to each of the *Soft-bubble* sensors.

in finding point-pair correspondences between the reference geometry and the sensed point-clouds. However, the non-smoothness of the SDF gradients on geometries with sharp edges or corners can prove detrimental to the optimization.

Therefore, we extend the notion of the SDF to a more generic scalar *proximity* field that exhibits smooth gradients outside a given geometry even when sharp corners or edges are present. This can be defined as the function ϕ , a scalar field that is specified for a given geometry with an associated reference frame G that maps any point P that is specified in a frame C to a scalar value c as in,

$$\phi({}^G X^C p^P) = c, {}^C p^P \in \mathbb{R}^3, \quad c \in \mathbb{R}^1, \quad (1)$$

where the notation ${}^C p^P$ denotes a position vector p between C and P , and ${}^G X^C$ denotes the transform between C and G . By convention $c \geq 0$ outside the geometry, $c < 0$ inside and $c = 0$ at the boundary. It is desirable that the gradients for a proximity field are well-defined and smooth for all ${}^G p^P \in \mathbb{R}^3$. For smooth surfaces, similar to SDFs, the gradient satisfies the Eikonal equation $\nabla \phi = 1$. However, it is important that the smoothness is also guaranteed outside a geometry which can be ensured by considering the nearest corner distance along with the nearest surface distance as depicted in Fig. 7 for the cylinder case.

Then, given a point-cloud measurement ${}^T p^Q_{set_i} = [{}^T p^Q_0, \dots, {}^T p^Q_N]$ that is specified in the end effector frame T , the problem in its simplest form can be posed as a nonlinear optimization given by,

$$\theta^* = \underset{\theta}{\operatorname{argmin}} \sum_{i=0}^N [\phi({}^G X_T(\theta) {}^T p_i^Q)]^2, \quad \text{s.t. : } \theta \in \text{SE}(3), \quad (2)$$

where θ is the chosen parameterization of the pose ${}^G X_T$, for a 7-dimensional vector composed of the three translational coordinates and four rotational coordinates (represented by

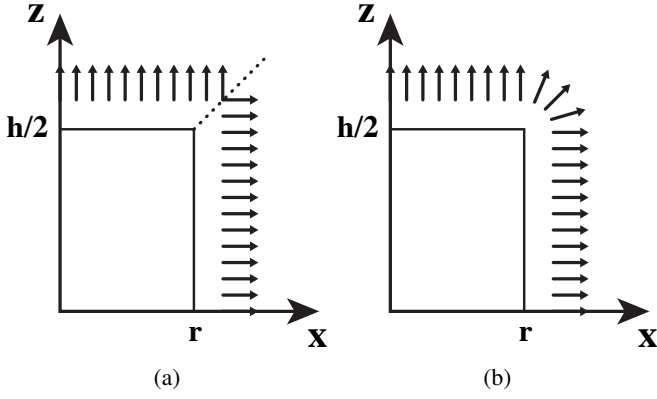


Fig. 7: Proximity field gradients illustrated on a 2D projection of a cylinder geometry; (a) signed-distance fields without corner correction and (b) proximity fields with corner correction.

the quaternion). We use a unit-norm constraint on the the quaternion part of θ .

Here, we specifically focus on analytic *proximity* fields for cylinders which we develop to ensure smoothness of $\nabla\phi$ outside the geometry. This function can be defined in the following manner:

Given a cylindrical geometry of height h and radius r , with the height along the z -axis, and the half-height at the origin of G , consider a slab δ_{sl} consisting of the planar projection of the cylinder, then the *proximity* field ϕ can be defined as,

$$\phi(Gp^P) = \|\max(\delta_{sl}, 0_2)\|_2 + v\max(\min(\delta_{sl}, 0_2)), \quad (3)$$

where the cylinder slab function δ_{sl} is defined by,

$$\begin{aligned} \delta_{sl} &= |d| - b, \\ d &= [\|Gp^P(x, y)\|_2, Gp^P(z)]^T, \\ b &= [r, h/2]^T, \end{aligned}$$

where the operators \max and \min are the binary operators returning the maximum and minimum of a pair of scalar values and $v\max$ denotes the operator that return the \max value among the components of a given vector, x, y , and z denote the components of Gp^P and the modulus operation $|d|$ is element-wise. Fig. 7b illustrates the field at some points outside the cylinder.

This algorithm possesses several advantages from the perspective of implementation and efficiency. First, the analytical form of the *proximity* fields for several basic primitive shapes can be easily derived and thus efficiently computed. Second, similar to SDFs, multiple primitive *proximity* fields can be aggregated for complex geometries and the gradients can be analytically computed. We use automatic differentiation to efficiently compute the gradients, and solve the constrained nonlinear optimization in Eq. 2 using SNOPT [20].

For *Soft-bubbles* we compute a concatenated grasp point-cloud $Tp^Q_{set_i}$ from the two individual point-clouds by using gripper dimensions and the joint-encoder measurements from the parallel gripper. In practice, the speed and accuracy of

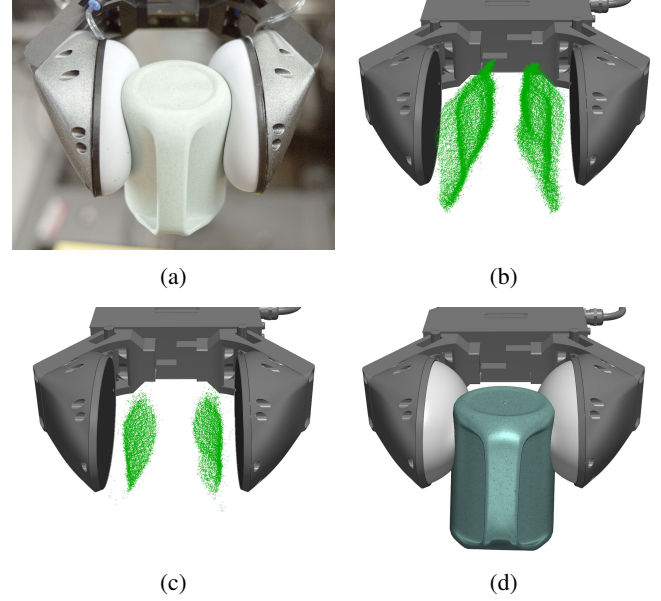


Fig. 8: The various stages of the in-hand pose estimation pipeline. (a) A plastic mug being grasped in the *Soft-bubble* gripper system. (b) The concatenated point-cloud produced from the depth images from each *Soft-bubble* sensor computed in the gripper frame. (c) The contact-patch filtered concatenated point-cloud (d) The estimated in-hand mug pose from the *proximity* pose estimator.

this *proximity* pose optimization were greatly enhanced by cropping this point-cloud to only include the contact-patch (similar to our earlier work [3]). For the results shown in this paper, since we only test shapes that are convex, we used a more naive contact patch filter that utilized a difference image generated with respect to a reference configuration of the bubbles when they are free from contact. The various stages of this entire pipeline can be seen in Fig. 8. In practice, the contact patch point-cloud size is typically around 10K points; the pose-estimation computation frequency tends to be around $0.5Hz - 1Hz$ on a standard multicore processor system using parallelization of computation of Eq. 2.

V. SHEAR DEFORMATION ESTIMATION

Incorporating a pseudorandom dot pattern on the interior of the bubble surface (Fig. 5) allows us to extract tangential displacement in the plane orthogonal to the applied gripper force ($P_{membrane}$).

As shown in Sec. III, the image sensor providing IR images is angled to look at the membrane. We extract a vector field using dense optical flow via the Gunner Farnebacks algorithm [21]. Since tracking algorithms can have difficulties handling sudden changes, we instead focus on the *relative* shear-displacement problem of comparing two images: a reference $I_{t=0}$ and a later image $I_{t=n}$. The reference image is taken after the gripper has achieved a stable grasp of the target object. The second $I_{t=n}$ image is the latest from the camera. The output is vector \vec{V} for each pixel of the image that indicates an

estimate of the displacement of that patch in pixels. To extract the tangential displacement vector \vec{s} , we sum the vector field on the image by $\vec{s} = \sum_{i=0, \vec{v}_i \in \vec{V}}^n \vec{v}_i$, where \vec{s} is the tangential displacement, \vec{v}_i is the vector value at pixel i , n is the pixel number of input image. The output relative shear displacement is then represented by the normalized form $\vec{s}/\|\vec{s}\|$.

VI. TACTILE OBJECT CLASSIFICATION

As described above, the ToF sensors in the bubbles provide both depth and IR over independent channels. We implemented a classifier that can utilize either of these streams from both cameras. Each camera’s resolution is 224×176 . Only stable grasps are used for training and inference - we compute if a grasp is stable using thresholds on the bubble pressure differential as well as the finger velocity. This eliminates confounding issues like image blur. The images from both bubbles are then concatenated (as in Fig. 9) and down-sampled into a 224×224 image. ResNet18 was used as the network architecture and both training and inference pipelines were implemented with *pytorch*.

A. Automatic Labelling Pipeline

In contrast to our previous classifier [2], we observed that when using images from two *Soft-bubbles* to classify objects during manipulation, it is important to collect training images of a sufficient diversity of valid stable grasps per class. Further, any unique geometric features of an object (e.g. handle of a mug, cap of a bottle) need to be sufficiently represented in the training data. Initial training attempts using human-generated grasp configurations resulted in insufficient samples of object geometries in poses the robot might grasp, resulting in poor classifier performance when tested on the robot. To mitigate these issues, we built an automatic labeling pipeline. One or more objects were placed into the robot’s workspace and the robot was commanded to repeatedly pick up objects using the same grasp generator that is used in our object sorting control pipeline (antipodal grasps on the object computed using an external depth sensor [22]) and then drop them randomly. Once a stable grasp is obtained, collecting about one minute of data (600 training samples) for each object was found to be sufficient for training.

This pipeline can generate 50 grasps (30K training samples) autonomously within 1.5 hours. We trained on three classes of objects (see Fig. 10). The training converged in about 45 epochs using stochastic gradient descent with a validation accuracy of $\geq 99\%$ (validated on a separate random samples of 30% of the labelled data). For both training and inference we utilized a dual GPU Intel Xeon workstation. More rigorous analysis of the number of grasps required are ongoing and out of the scope of this paper.

VII. EXPERIMENTS

To demonstrate robust manipulation we deployed the *Soft-bubble* gripper and the presented tactile perception techniques in a cluttered indoor scenario. The experimental setup consisted of a fixed-base manipulator in a realistic kitchen setup

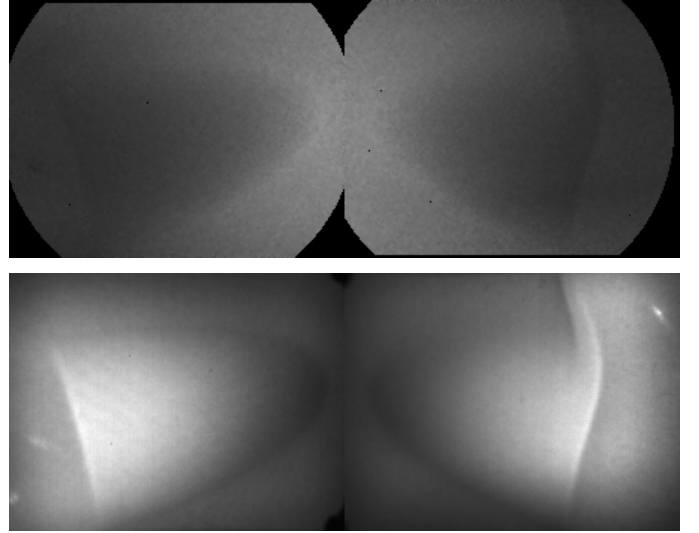


Fig. 9: An instance of an automatically labelled training sample for the plastic mug; Top: the concatenated depth image; Bottom: the concatenated IR image.

with counters, a sink, and a dishwasher. The setup and the various tasks that we demonstrated can be seen in this [video](#).

A. Robot and Task Setup

The bubble fingers were mounted on a parallel Schunk WSG 50 gripper attached to a KUKA iiwa arm. A set of six Intel RealSense D415 cameras were employed for external (visual) perception of the sink-scene; three cameras pointed into the sink at a downward angle, an additional two cameras facing the counter-top and one final camera tracking the three racks (using fiducial markers) and the door of the dishwasher. The following objects were used: plastic mugs, plastic PET bottles (acting as refuse in the sink), and wine glasses. An OptiTrack motion capture system was used for validating the in-hand pose estimation.

Depth images from the RealSense cameras were fused to generate a combined point cloud of the scene within the sink. Grasp poses were identified on objects either by an anti-podal grasp generation pipeline [22] or by a naive blob centroid location.

Motion planning was implemented using a combination of sampling-based and optimization-based techniques [23]. The robot’s actions were scripted in the form of a set of parameterized closed-loop *primitives*. Real-time control ran over Ethernet/IP using KUKA’s FRI interface, and the software framework used a mixture of tools written in C++ and Python. The Drake simulation and control toolbox [24] was used for the kinematic computations within the control and planning stack. The collision detection in planning and control used the FCL collision library (through Drake). All of the perception, planning, and control software was running on one Intel Xeon workstation under Ubuntu 18.04.



Fig. 10: Some of the grasped objects from the dish loading task and their corresponding *Soft-bubble* IR images. In each of these cases, the object class was correctly identified with a probability $p > 0.95$; (a) plastic PET tea-bottle, (b) square form-factor plastic alcohol bottle, (c) plastic mug, and (d) an alternative grasp utilized on a bottle similar to that in (b).

B. Task 1: Dish loading with recycling separation

In this task the robot autonomously reaches into the sink, grasps objects, and classifies them as either a valid dish (e.g. mugs) or recycling (e.g. PET bottles). The limit on the objects selected was due to our existing dish-loading pipeline, not fundamental limitations of *Soft-bubbles*. Valid dishes must be placed within the dishwasher and recycling dumped in a bin. screen

The images seen in Fig. 10 depict some of the different classes of objects that we tested on the dish-loading setup during the post-grasp tactile classification stage.

C. Task 2: Pose Estimation Under Antagonism

We analyzed the performance of the *proximity* pose optimizer for cylindrical mugs by equipping the manipulant with active motion capture markers. For an arbitrarily chosen in-hand pose, the pose estimates computed were compared against the motion capture output. Since the *proximity* field we use here is that of a cylinder, the pose-estimates can be considered reliable only for contact with cylindrical surfaces. The peak range of the offset in the seed from the ground-truth pose, which still exhibited stable convergence, was observed to be ± 30 degrees in pitch and roll directions. We used either the pre-grasp pose estimated by our vision system (when available) or seeded it with a cardinal pose. Once a stable grasp was achieved, we initialized our pose optimization with the previous estimated pose for continuous tracking.

D. Task 3: Shear-based Manipulant Release and Handover

We implemented a closed-loop controller that monitors the shear forces on a manipulant relative to an initial stable grasp. By comparing it to a threshold we implemented automated

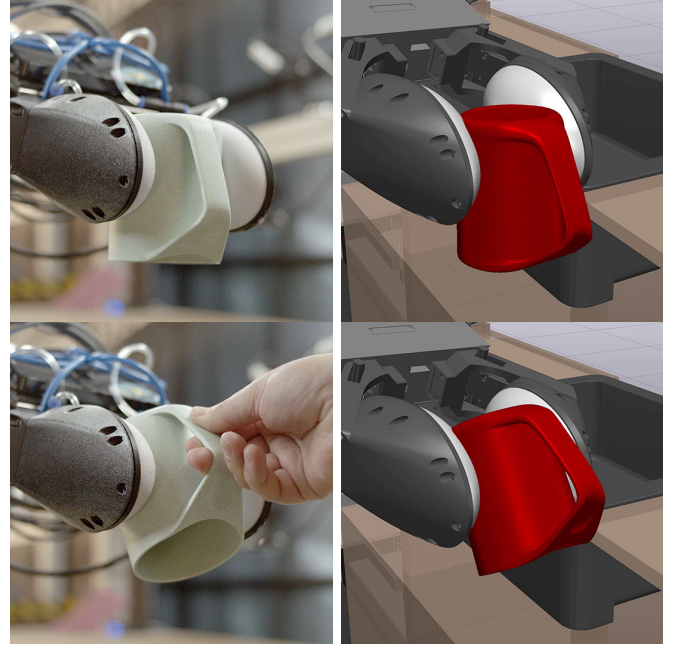


Fig. 11: Mugs grasped in two different poses and their corresponding in-hand pose-estimates; Top: An upside down mug; Bottom: Grasped mug being perturbed. The pose estimator runs at 1-2 Hz and can track these kinds of perturbations.

placing or handover of the manipulant. This was demonstrated with two tasks. In the first we hand a fragile object (wine glass with liquid) to a human as seen in Fig. 12. The second is an autonomous blind stacking of wine glasses where the robot has no notion of the geometry of the objects or the stacking surface. This latter task can be seen in this [video link](#).

The thresholds used are dependent on the post-grasp contact patch size and geometry. Future work will involve using the contact patch estimator and pressure sensor data to improve shear sensing accuracy for more diverse objects.

VIII. CONCLUSIONS AND DISCUSSION

In this paper, we present the *Soft-bubble* gripper system, which combines a highly compliant gripping surface and visuotactile sensing in the form factor of a parallel gripper. The tactile sensing capability of the gripper enables multiple forms of perception; in multiple real-world tasks, we demonstrated in-hand pose-estimation, shear-deformation estimation and tactile classification, within a robust manipulation pipeline.

As seen in our results, the combination of *Soft-bubbles* and tactile perception can be valuable for solving manipulation problems in cluttered environments. The *Soft-bubble* fingers are affordable, compliant, easy to construct and mechanically robust. In addition, the perception methods presented are computationally efficient enough to enable closed-loop, real-time control of complex tasks. The *proximity* field method used in our pose estimation framework is a novel contribution that is promising for achieving tractable depth-based tracking.

Soft-bubbles still have room for improvement. Our current implementation is limited in size due to the minimum range

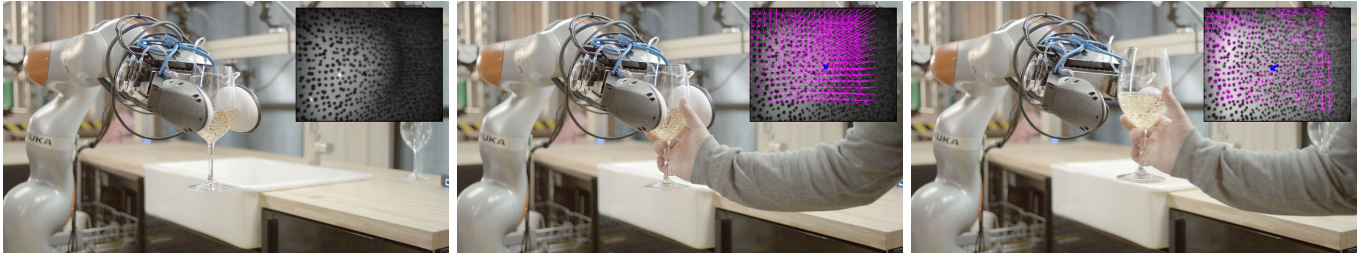


Fig. 12: Video frames from a shear-based manipuland handover task with a liquid filled fragile manipuland (wine-glass). The shear displacement estimator’s display output (computed from one of the *Soft-bubble* fingers) is overlaid on the top right of each panel: (left) Before handover, showing the nominal configuration. From this point on, shear is computed relative to this configuration. (center) At the onset of handover, when a human exerts forces to take the manipuland from the robot, relative shear displacement increases until a threshold on the norm of the magnitude. (right) Condition immediately after release; there is residual relative shear seen with respect to the pre-handover condition.

of custom ToF sensors that are difficult to source, so we are working on customized ToF cameras designed specifically for short-range depth sensing.

Algorithmically, we are working on enhancements to both shear displacement estimation and pose tracking. Our current shear displacement tracking is based on tracking visual features with the IR channel. By incorporating the depth channel from the ToF camera we hope to have the *Soft-bubble* act more like a full low-cost force-torque sensor. Our pose estimator is being expanded to handle other geometric primitives (e.g. frustrums).

Lastly, we are investigating more challenging applications for the *Soft-bubble* technology, including tool use and enabling physical interaction between humans and robots.

REFERENCES

- [1] J. Hughes, U. Culha, F. Giardina, F. Guenther, A. Rosendo, and F. Iida, “Soft manipulators and grippers: A review,” *Frontiers in Robotics and AI*, vol. 3, p. 69, 2016. [Online]. Available: <https://www.frontiersin.org/article/10.3389/frobt.2016.00069>
- [2] A. Alspach, K. Hashimoto, N. Kuppawamy, and R. Tedrake, “Soft-bubble: A highly compliant dense geometry tactile sensor for robot manipulation,” in *2019 2nd IEEE International Conference on Soft Robotics (RoboSoft)*. IEEE, 2019, pp. 597–604.
- [3] N. Kuppawamy, A. Castro, C. Phillips-Grafflin, A. Alspach, and R. Tedrake, “Fast model-based contact patch and pose estimation for highly deformable dense-geometry tactile sensors,” *IEEE Robotics and Automation Letters*, 2019.
- [4] K. Shimonomura, “Tactile image sensors employing camera: A review,” *Sensors*, vol. 19, no. 18, p. 3933, 2019.
- [5] J. Kim, A. Alspach, and K. Yamane, “3d printed soft skin for safe human-robot interaction,” in *IROS*. IEEE, 2015, pp. 2419–2425.
- [6] W. Yuan, S. Dong, and E. H. Adelson, “GelSight: High-resolution robot tactile sensors for estimating geometry and force,” *Sensors*, vol. 17, no. 12, p. 2762, 2017.
- [7] A. Wilson, S. Wang, B. Romero, and E. Adelson, “Design of a Fully Actuated Robotic Hand With Multiple Gelsight Tactile Sensors,” *arXiv:2002.02474 [cs]*, Feb. 2020, arXiv: 2002.02474. [Online]. Available: <http://arxiv.org/abs/2002.02474>
- [8] E. Donlon, S. Dong, M. Liu, J. Li, E. Adelson, and A. Rodriguez, “GelSlim: A high-resolution, compact, robust, and calibrated tactile-sensing finger,” *arXiv preprint arXiv:1803.00628*, 2018.
- [9] W. Yuan, Y. Mo, S. Wang, and E. H. Adelson, “Active clothing material perception using tactile sensing and deep learning,” in *ICRA*, May 2018, pp. 1–8.
- [10] Y. Zhang, Z. Kan, Y. Yang, Y. A. Tse, and M. Y. Wang, “Effective estimation of contact force and torque for vision-based tactile sensor with helmholtz-hodge decomposition,” *CoRR*, vol. abs/1906.09460, 2019. [Online]. Available: <http://arxiv.org/abs/1906.09460>
- [11] J. Issac, M. Wüthrich, C. G. Cifuentes, J. Bohg, S. Trimpe, and S. Schaal, “Depth-based object tracking using a robust gaussian filter,” in *2016 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2016, pp. 608–615.
- [12] M. C. Koval, M. R. Dogar, N. S. Pollard, and S. S. Srinivasa, “Pose estimation for contact manipulation with manifold particle filters,” in *IROS*. IEEE, 2013, pp. 4541–4548.
- [13] T. Schmidt, R. A. Newcombe, and D. Fox, “Dart: Dense articulated real-time tracking,” in *Robotics: Science and Systems*, vol. 2, no. 1. Berkeley, CA, 2014.
- [14] T. Schmidt, K. Hertkorn, R. Newcombe, Z. Marton, M. Suppa, and D. Fox, “Depth-based tracking with physical constraints for robot manipulation,” in *2015 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2015, pp. 119–126.
- [15] R. Li, R. Platt, W. Yuan, A. ten Pas, N. Roscup, M. A. Srinivasan, and E. Adelson, “Localization and manipulation of small parts using gelsight tactile sensing,” in *2014 IEEE/RSJ International Conference on Intelligent Robots and Systems*. IEEE, 2014, pp. 3988–3993.
- [16] M. Bauza, O. Canal, and A. Rodriguez, “Tactile mapping and localization from high-resolution tactile imprints,” *arXiv preprint arXiv:1904.10944*, 2019.
- [17] G. Izatt, G. Mirano, E. Adelson, and R. Tedrake, “Tracking objects with point clouds from vision and touch,” in *ICRA*. IEEE, 2017, pp. 4000–4007.
- [18] H. Liu, Y. Wu, F. Sun, and G. Di, “Recent progress on tactile object recognition,” *International Journal of Advanced Robotic Systems*, vol. 14, pp. 1–12, 07 2017.
- [19] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *CVPR*, 2016, pp. 770–778.
- [20] P. E. Gill, W. Murray, and M. A. Saunders, “Snopt: An sqp algorithm for large-scale constrained optimization,” *SIAM review*, vol. 47, no. 1, pp. 99–131, 2005.
- [21] G. Farneback, “Two-frame motion estimation based on polynomial expansion,” in *SCIA*, 2003.
- [22] A. ten Pas, M. Gualtieri, K. Saenko, and R. Platt, “Grasp pose detection in point clouds,” *The International Journal of Robotics Research*, vol. 36, no. 13-14, pp. 1455–1473, 2017.
- [23] M. Fallon, S. Kuindersma, S. Karumanchi, M. Antone, T. Schneider, H. Dai, C. P. D’Arpino, R. Deits, M. DiCicco, D. Fourie, T. Koolen, P. Marion, M. Posa, A. Valenzuela, K.-T. Yu, J. Shah, K. Iagnemma, R. Tedrake, and S. Teller, “An architecture for online affordance-based perception and whole-body planning,” *Journal of Field Robotics*, vol. 32, no. 2, pp. 229–254, 2015.
- [24] R. Tedrake and the Drake Development Team, “Drake: Model-based design and verification for robotics,” 2019. [Online]. Available: <https://drake.mit.edu>