

MidasTouch: Monte-Carlo inference over distributions across sliding touch

Sudharshan Suresh^{1,2}, Zilin Si¹, Stuart Anderson², Michael Kaess¹, Mustafa Mukadam²

¹Carnegie Mellon University, ²Meta AI

<https://suddhu.github.io/midastouch-tactile>

Abstract: We present MidasTouch, a tactile perception system for online global localization of a vision-based touch sensor sliding on an object surface. This framework takes in posed tactile images over time, and outputs an evolving distribution of sensor pose on the object’s surface, without the need for visual priors. Our key insight is to estimate local surface geometry with tactile sensing, learn a compact representation for it, and disambiguate these signals over a long time horizon. The backbone of MidasTouch is a Monte-Carlo particle filter, with a measurement model based on a tactile code network learned from tactile simulation. This network, inspired by LIDAR place recognition, compactly summarizes local surface geometries. These generated codes are efficiently compared against a pre-computed tactile codebook per-object, to update the pose distribution. We further release the YCB-Slide dataset of real-world and simulated forceful sliding interactions between a vision-based tactile sensor and standard YCB objects. While single-touch localization can be inherently ambiguous, we can quickly localize our sensor by traversing salient surface geometries.

Keywords: Tactile perception, Localization, 3D deep learning

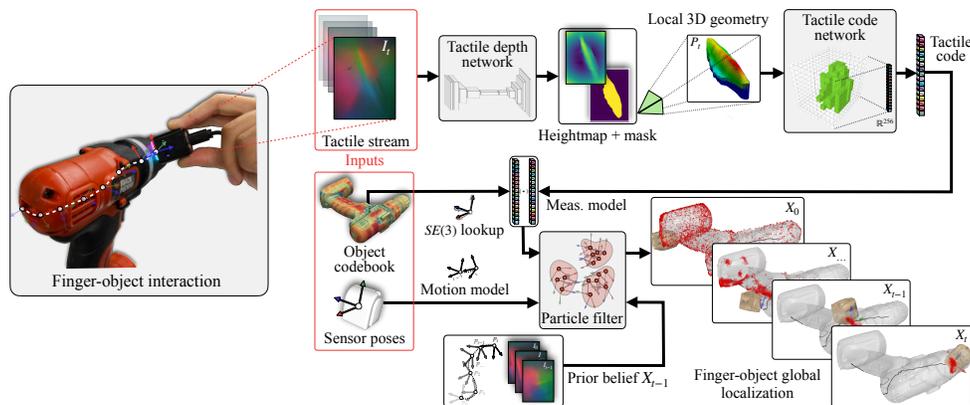


Figure 1: MidasTouch performs online global localization of a vision-based touch sensor on an object surface during sliding interactions. Given posed tactile images over time, this system leverages local surface geometry within a nonparametric particle filter to generate an evolving distribution of sensor pose on the object’s surface.

1 Introduction

Interactive perception is the Catch-22 [1] of robotics; vision can be used to track objects for downstream contact-rich interactions, but such interactions can occlude visual tracking. In particular, end-effector-object relative positioning is crucial for contact-rich policies like sliding [2–6], in-hand manoeuvres [7–10], finger-gaiting [11], and multi-finger enclosure [12]. Additionally, vision is often affected by object transparency, specularities, and poor scene illumination. With high-dimensional vision-based tactile sensors [13–15], we now have a window into local object interactions. While tactile images from these sensors capture local surface geometry, they lack global context necessary for relative pose tracking. We address this challenge with **MidasTouch**, an online tactile perception system that tracks the evolving pose distribution of a vision-based touch sensor sliding across known objects. We acquire global context by integrating local observations over long horizons.

The form-factor of vision-based tactile sensors has restricted prior methods to small parts [16–18] or local tracking [19, 20]. For everyday objects, we can leverage compositionality [21]; each is made up of local geometries that are spatially fixed relative to one another. Consider a simple mug: made up of a curved body, flat base, rounded handle, and sharp lip. Without global context, a single-touch is ambiguous: a perceived sharp edge could lie anywhere along the lip of the mug. Such a likelihood distribution is spread across the object’s surface and not unimodal, but interaction over long time horizons can disambiguate it. This mirrors haptic *apprehension*, or the exploratory procedures humans perform when presented with a familiar object [22].

We approach this as an analogue to mobile robot Monte-Carlo filtering [23], but instead apply it on the surface manifold of the object. Just as a mobile robot has access to detailed floor plans, odometry, and cameras, manipulators have access to object meshes, end-effector poses, and vision-based touch. This is applicable to environments with known object models like households, warehouses, and factories, further facilitated by large-scale scanned object datasets [24, 25]. While priors from vision can bootstrap our system [26], they are not a prerequisite for global estimation.

In addition to open-sourcing MidasTouch, we release a comprehensive real-world and simulated dataset of sliding DIGIT [15] interactions across standard YCB objects [24] with ground-truth. This serves as both an evaluation of MidasTouch, and as a benchmark currently lacking in the tactile sensing community. Both are accessed from our project [website](#). Specifically, our contributions are:

1. Online particle filtering over posed tactile images for a distribution of finger-object poses,
2. Learned embeddings for vision-based touch using local surface geometry,
3. The YCB-Slide dataset of DIGIT and YCB object interactions for evaluation and benchmarking.

Our framework relies on recent developments in tactile sensing, rendering, and robotics: **(i)** vision-based tactile sensors [13–15, 27–31] like the GelSight and DIGIT have the requisite spatial acuity to discern local geometric features, **(ii)** tactile simulation [32–34] with realistic rendering enables learning tactile observation models and precomputing object-specific interactions with sizeable data, which would be infeasible in the real-world, and **(iii)** learned models for 3D place recognition [35–37] spawned by ubiquitous LIDAR and RGB-D data that can be extended to tactile sensing.

2 Related work

Tactile pose inference: Binary contact sensors have been used in conjunction with particle filters for global estimation of robot hands relative to simple geometries [38–43]. These methods require a large amount of touches, but serve as a touchstone for Monte-Carlo methods that have seen great success in mobile robotics [23, 44]. With tactile arrays, local patch measurements are integrated for planar estimation [45, 46] or for 3D alignment [47]. The synergy of vision and touch gives much needed global context [48, 49] but is outside our current work’s scope.

With vision-based touch, Li et al. [16] show planar small-part localization by directly computing the homography on GelSight heightmaps. Relative pose-tracking has also been learnt directly from tactile images either via recurrent [50], or auto-encoder [15, 19] networks. Additionally, local tracking has been explored with online factor-graph optimization [19, 20, 51]. The aforementioned approaches are unimodal and rely on good pose initializations; MidasTouch is nonparametric and requires no such knowledge. Additionally, Kelestemur et al. [52] show category-level object pose estimation with a parallel jaw gripper.

Closely related is Tac2Pose [17, 18], which performs pose estimation of small-parts with the GelSlim [14]. It produces a distribution of object poses learned in tactile simulation from contact shapes. Similarly, Gao et al. [53] estimates contact location on objects through vision-based touch and audio with a small, discrete set of measurements. We differentiate our work in both context and approach: **(i)** we filter over long time horizons rather than single/multi touch predictions, **(ii)** tactile embeddings are learned from local surface geometry rather than images, and **(iii)** we consider objects considerably larger than the robot finger.

Place recognition for touch: A compact representation for tactile images enables easy frame-to-frame or frame-to-model tracking. Inspired by the computer vision community, a popular interme-

diate is a learned embedding from either RGB [15, 19] or binary contact masks [17, 18]. Realistic tactile simulators, like TACTO [32], enable training these models to the scale and generality of arbitrary real-world interactions. For example, Bauza et al. [18] learn object-specific embeddings to match 2D contact shapes.

For contact-rich manipulation, local 3D geometry is a natural candidate, and iterative closest point (ICP) has shown promise for frame-to-frame tracking [20, 54]. For vision-based touch, 3D geometry is obtained either via photometric stereo [13, 55, 56], or image-to-heightmap models [20, 31, 57–59]. However, ICP is only suitable for local tracking and not global localization: it is sensitive to initialization and is intractable to scale as a sampling-based measurement model.

Succinctly, what is an accurate and efficient similarity metric for tactile geometry? With the ubiquity of LIDAR/RGB-D data, matching unordered point sets is vital for place recognition in the SLAM community [36, 37]. Following the seminal PointNet [60], Choy et al. [61] later developed efficient and expressive sparse 3D convolution. This backbone was used to aggregate local point features to global point-cloud embeddings for LIDAR place-recognition [37]. We show tactile geometries can be compacted into codes with the same architecture, for easy frame-to-model queries.

3 Problem formulation

For a vision-based tactile sensor dynamically sliding along a known object’s surface, our goal is to track the distribution of 6D sensor pose $\mathbf{x}_t \in SE(3)$ in the object-centric frame. Such pose context is useful for downstream planning and control, for instance manipulating an object in-hand [7]. The sensor is affixed on a robot finger and we have access to the end-effector pose, but its relative position and orientation on the object’s surface is unknown. At each timestep t , our measurements are a tactile image, and noisy sensor pose in the robot’s reference-frame, $\mathbf{z}_t = \{\mathbf{I}_t \in \mathbb{R}^{240 \times 320 \times 3}, \mathbf{p}_t \in SE(3)\}$. For simplicity, we assume the object is stationary and sliding is a single continuous contact interaction. The dimensions of the objects are much larger than the sensor’s contact area. We assume an uninformative initial prior for \mathbf{x}_t , but later show the benefit of visual priors.

4 Global localization during sliding touch

MidasTouch comprises three distinct modules, as illustrated in Figure 1: a tactile depth network (TDN: Section 4.1), tactile code network (TCN: Section 4.2), and particle filter (Section 4.3). At a high-level, the TDN first converts a tactile image \mathbf{I}_t into its local 3D geometry \mathbf{P}_t via a learned observation model. The 3D information is then condensed into a tactile code \mathbf{E}_t by the TCN through a sparse 3D convolution network. Finally, the downstream particle filter uses these learned codes in its measurement model, and outputs a sensor pose distribution that evolves over time. Throughout this work, we use 10 YCB objects with diverse geometries in our tests: sugar_box, tomato_soup_can, mustard_bottle, bleach_cleanser, mug, power_drill, scissors, adjustable_wrench, hammer, and baseball. In Appendix G, we further demonstrate the method on small parts.

4.1 Tactile depth network (TDN)

The tactile depth network learns the inverse sensor model to recover local 3D geometry from a tactile image. We adapt a fully-convolutional residual network [59, 62], and train it in a supervised fashion to predict local heightmaps from tactile images.

Network and training: This network uses a ResNet-50 backbone, with a series of up-projection blocks. With the optical tactile simulator TACTO [32], we render a large collection of DIGIT images with ground-truth heightmaps. The images are from simulated interaction with 40 YCB [24] object meshes, while holding out all test objects. The images are rendered at 5000 different poses per object: randomizing for contact point, orientation, and indentation depth. For sim2real transfer, we

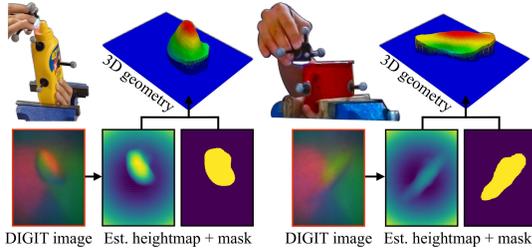


Figure 2: Real-world DIGIT images from interactions in the YCB-Slide dataset (more in Appendix A). Our tactile depth network (Section 4.1) is trained in simulation, and predicts 3D geometry given input tactile image.

calibrate TACTO with real-world data from multiple independent DIGITs, and apply random lighting augmentations (refer Appendix A). Per object, the split of train-validation-test is 4000:500:500.

Image to 3D: At each timestep, we pass the 240×320 RGB tactile image \mathbf{I}_t through the network to get heightmap \mathbf{H}_t . To remove non-contact areas we compute a mask \mathbf{C}_t by depth thresholding, and get $\hat{\mathbf{H}}_t = \mathbf{H}_t \odot \mathbf{C}_t$. Finally, the heightmap is reprojected to 3D via the camera’s known perspective projection model $\hat{\mathbf{H}}_t \mapsto \mathbf{P}_t$. Evaluations on the test set show heightmap RMSE of 0.135 mm with respect to ground-truth. Figure 2 shows examples of the TDN on real-world DIGIT interactions.

4.2 Tactile code network (TCN)

This network summarizes large, unordered point clouds of local geometry into a low-dimensional embedding space, or code. If two sensor measurements are nearby in pose-space, they observe similar geometries and therefore, their codes will also be nearby in embedding-space. However, the inverse is not necessarily true: measurements from two opposite corners of a cube may have similar codes, but their corresponding sensor poses are dissimilar. This is an inherent challenge of tactile localization, the so-called *contact non-uniqueness* highlighted by Bauza et al. [18]. In our work getting the most-likely modes of this distribution is sufficient since the downstream particle filter can then disambiguate them temporally.

An analog in the SLAM community is the loop closure problem from 3D data. LIDAR place recognition modules, such as PointNetVLAD [36] and MinkLoc3D [37], use metric learning to learn point cloud similarity. While tactile images differ from natural images, the geometries from LIDAR and tactile sensors (normalized for scale) are similar. The MinkLoc3D architecture, based on MinkowskiNet [61], performs state-of-the-art point cloud retrieval through sparse 3D convolutions.

Network and training: The architecture comprises of three components: (i) Voxelization of \mathbf{P}_t into a quantized sparse tensor $\hat{\mathbf{P}}_t = \{\langle \hat{x}_i, \hat{y}_i, \hat{z}_i, 0/1 \rangle\}$, (ii) Feature pyramid network [63] for per-voxel local features $\hat{\mathbf{P}}_t^f = \{\langle \hat{x}_i, \hat{y}_i, \hat{z}_i, \hat{\mathbf{f}}_i \in \mathbb{R}^{256} \rangle\}$, and (iii) Generalized-mean pooling for point cloud embedding vector $\mathbf{E}_t \in \mathbb{R}^{256}$. We use the pre-trained weights, learned from four comprehensive LIDAR datasets [37], and fine-tune them with TACTO data collected over the 40 YCB objects. This is trained with a contrastive triplet loss, and we accumulate positive and negative local 3D geometries through pose-supervision.



Figure 3: Tactile codebook per object visualized as a spectral colorspace map using t-SNE [64]. Each codebook comprises of 50k densely sampled poses with their corresponding 256-dimensional tactile code. Similar hues denote sensor poses that elicit similar tactile codes. We can clearly delineate local geometric features: edges (sugar_box), ridges (power_drill), corners (scissors), and complex texture (baseball).

Tactile codebook: Once we have a code \mathbf{E}_t , we would like to efficiently compare this with a dense set of contacts. Inspired by prior work [18, 26], we build a tactile codebook comprising of $M = 50k$ randomized sensor poses on each object’s mesh. We evenly sample these points and normals on the mesh with random orientations and indentations (refer Appendix A). We feed these poses into TACTO to generate a dense set of tactile images for each object.

We pass the generated images through the TDN + TCN to get a codebook for the specific object o : $\mathcal{C}_o = \{\langle \mathbf{p}^{[m]}, \mathbf{E}^{[m]} \rangle\}_{m=1}^M$, where $\mathbf{p}^{[m]}$ are $SE(3)$ sensor poses in codebook and $\mathbf{E}^{[m]}$ are corresponding tactile codes. Figure 3 shows a t-SNE visualization of the codebooks, a colorspace representation of local geometric similarity. We observe regions with similar geometries have identical hues and the traversal of our sensor between these geometries will give us valuable measurement signals.

We build a KD-Tree with 6-element vectors: $\{\langle \mathbf{p}_{\text{trans}}^{[m]}, \alpha \log(\mathbf{p}_{\text{rot}}^{[m]}) \rangle\}_{m=1}^M$ for nearest-neighbor search. Here, $\log(\cdot)$ is the $SO(3)$ logarithm map obtained via Theseus [65], and $\alpha = 0.01$ is the rotation scaling factor. Thus, given a candidate pose $\mathbf{x}_t^{[i]}$ we do not have to render and obtain

its corresponding code, but instead just perform a pose-space lookup $\mathcal{C}_o(\mathbf{x}_t^{[i]})$. With memoization of code generation, we make getting codes for thousands of arbitrary pose particles tractable.

Single-touch localization: To understand the effectiveness of codes as a proxy for pose prediction, we conduct a set of single-touch experiments in simulation. This also provides an idea of which objects are salient, and which are adversarial, as highlighted in Figure 4. We observe that we perform significantly better than random touches for all objects, and those with symmetric, regular structures exhibit long-tail errors. This shows that while a single-touch has meaningful signal, it can be valuable to disambiguate these readings temporally with a filtering framework. We compare against embeddings from tactile images in Appendix B.

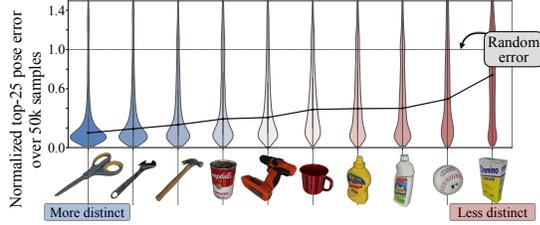


Figure 4: Pose-error for 50k single-touch queries on YCB objects, in ascending order (normalized with respect to random touch). For each query, we get the top-25 highest scores from the tactile codebook \mathcal{C}_o , and compute their minimum pose-error with respect to ground-truth. We observe tools with salient geometries to be easier to localize versus objects that exhibit symmetry.

4.3 Filtering over posed touch

The particle filter approximates the posterior distribution of the sensor pose as a set $\mathbf{X}_t = \{(\mathbf{x}_t^{[i]}, w_t^{[i]})\}_{i=1}^{N_t}$, where $\mathbf{x}_t^{[i]}$ are $SE(3)$ sensor pose particles, $w_t^{[i]}$ are the predicted particle weights, and N_t is number of particles. Rather than a unimodal Gaussian, this arbitrary distribution captures the multi-modality of global localization. We propagate the posterior distribution based on the stream of tactile images and noisy sensor poses, emulating a conventional particle filter but with a learned measurement model. Additional implementation details can be found in Appendix C.

Initialization: The initial distribution \mathbf{X}_0 is sampled from a coarse prior, which can (optionally) be obtained from vision [26]. Estimating contact location from vision is inherently noisy and can be ambiguous due to object symmetries. In our experiments, we find that even with very uninformative priors, our multi-modal filter can prune out hypotheses and converge to the correct one.

We sample particles from a 6D prior about the ground-truth pose \mathbf{x}_0^{gt} as $\mathbf{x}_0 \sim \mathcal{N}(\mathbf{x}_0^{\text{gt}}, \begin{bmatrix} \sigma_{\text{trans}} & 0 \\ 0 & \sigma_{\text{rot}} \end{bmatrix})$, such that $3\sigma_{\text{trans}} = \mathcal{M}_{\text{diag}}$ and $3\sigma_{\text{rot}} = 180^\circ$. This weak prior scatters the particles across the object surface, and allows comparison across objects as a function of their mesh diagonal length $\mathcal{M}_{\text{diag}}$. In our results, we also show particle filter ablations with tighter pose priors. After sampling, we project the particles back onto the surface by querying their nearest neighbors from the tactile codebook.

Motion model: We typically have access to noisy global estimates of end-effector pose, p_t , via robot kinematics. Our sensor odometry comes from the relative pose predictions $\Delta p_t = (p_{t-1})^{-1} \cdot p_t$. The motion model propagates particles \mathbf{X}_{t-1} forward by sampling from a state transition probability:

$$\mathbf{x}_t^{[i]} \sim p(\mathbf{x}_t | \mathbf{x}_{t-1}^{[i]}, \Delta p_t) = \mathbf{x}_{t-1}^{[i]} \oplus \mathcal{N}(\Delta p_t, \Sigma_{\Delta p}) \quad (1)$$

Additionally, odometry can also be from marker flow [66] or image-to-image tracking [19, 20].

Measurement update: In this step, we update the particle weights based on how their tactile codes match the current measurement \mathbf{E}_t , as a function of cosine distance. At each timestep, we lookup the codebook for particle codes (Section 4.2) and perform a matrix-vector multiplication:

$$w_t^{[i]} = p(\mathbf{z}_t | \mathbf{x}_t^{[i]}) \sim \text{softmax} \left(\frac{\mathbf{E}_t \cdot \mathcal{C}_o(\mathbf{x}_t^{[i]})}{\|\mathbf{E}_t\| \cdot \|\mathcal{C}_o(\mathbf{x}_t^{[i]})\|} \right) \quad (2)$$

These weights are scaled $[0, 1]$ and represent how well the candidate particle poses match with the current measurement. These weights feed into the subsequent resampling step.

Particle resampling: We sample a new set of particles \mathbf{X}_t from the proposal distribution with a probability proportional to their weights. Low-variance resampling [67] is a popular solution that covers the sample set in a systematic manner.

Hypothesis clustering: Alongside a distribution of poses, downstream tasks may need distinct pose hypotheses. To achieve this, we hierarchically cluster particles in \mathbb{R}^3 using DB-SCAN [68]. We average the cluster positions and quaternions [69] to get a hypothesis set $h_t : \{\mathbf{x}_t^1 \dots \mathbf{x}_t^H\}$.

5 The YCB-Slide dataset

To evaluate MidasTouch, and enable further research in tactile sensing, we introduce our YCB-Slide dataset (more in Appendix D). It comprises of DIGIT sliding interactions on the 10 YCB objects from Section 4. We envision this can contribute towards efforts in tactile localization, mapping, object understanding, and learning dynamics models. We provide access to DIGIT images, sensor poses, ground-truth mesh models, and ground-truth heightmaps + contact masks (simulation only).

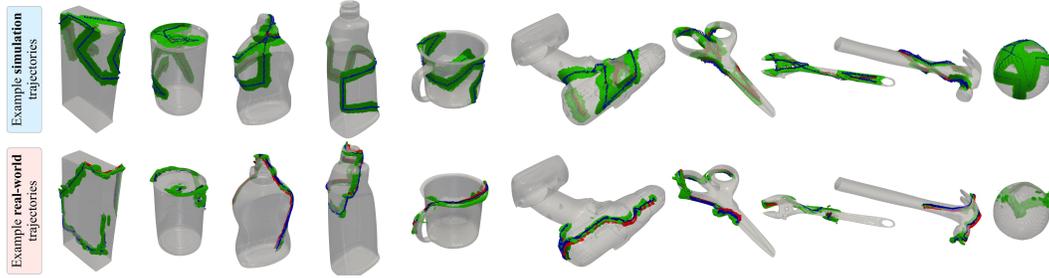


Figure 5: Example sliding trajectories from simulated and real trials on the 10 YCB objects. Overlaid in green are the local 3D geometries captured by the tactile sensor, and the contact poses as RGB coordinate axes.

Simulated interactions: We simulate sliding across objects using TACTO, mimicking realistic interaction sequences. The pose sequences are geodesic-paths of fixed length $L = 0.5\text{m}$, connecting random waypoints on the mesh. We corrupt sensor poses with zero-mean Gaussian noise $\sigma_{\text{trans}} = 0.5\text{mm}$, $\sigma_{\text{rot}} = 1^\circ$. We record five trajectories per-object; 50 interactions in total. A representative sample of the trajectories, along with their accrued tactile geometries are shown in Figure 5.

Real-world interactions: In the real-world, we perform sliding experiments through handheld operation of the DIGIT. We keep each YCB object stationary with a heavy-duty bench vise, and slide along the surface and record 60s trajectories at 30Hz. We use an OptiTrack system for timesynced sensor poses, with 8 cameras tracking the reflective markers. We affix six markers on the DIGIT and four on each test object. The canonical object pose is adjusted to agree with the ground-truth mesh models. Minor misalignment of sensor poses from human error are rectified by postprocessing the trajectories to lie on the object surface. We record five logs per-object, for a total of 50 sliding interactions. Our experimental setup and dataset are illustrated in Figures 5 and 6.



Figure 6: Real-world sliding trials in the YCB-Slide dataset. Inset is an example tactile image from the interactions, capturing the local geometry of the object.

6 Experimental results

Simulation: We evaluate MidasTouch over the 50 trajectories collected in Section 5. As the filter is non-deterministic, each trajectory is run 10 times for averaged results over 500 trials. Figure 7 shows qualitative results for three representative trajectories. At each timestep, we visualize the pose distribution and plot the average particle RMSE with respect to the ground-truth pose. We see convergence to the most-likely pose hypothesis over the sliding interactions. Upon convergence, the hypothesis clustering step sets an averaged pose for each cluster: visualized as the pose axes with uncertainty ellipses. We also highlight the current tactile image, surface geometry, and comparison with the tactile codebook. These trials are initialized with pose prior $[\sigma_{\text{trans}}, \sigma_{\text{rot}}]$ from Section 4.3.

Figure 8 [left] accumulates quantitative metrics over the 500 simulated trials. First, we show the final pose errors across all trials compared against the initial error. Overall, the averaged final pose errors are 0.74cm and 9.43° ; the per-object errors roughly correlate to the single-touch errors from Figure 4. While the error drops significantly for most trials, it fluctuates depending on each trajectory's salient geometries (or lack thereof). Larger pose-errors can be at-

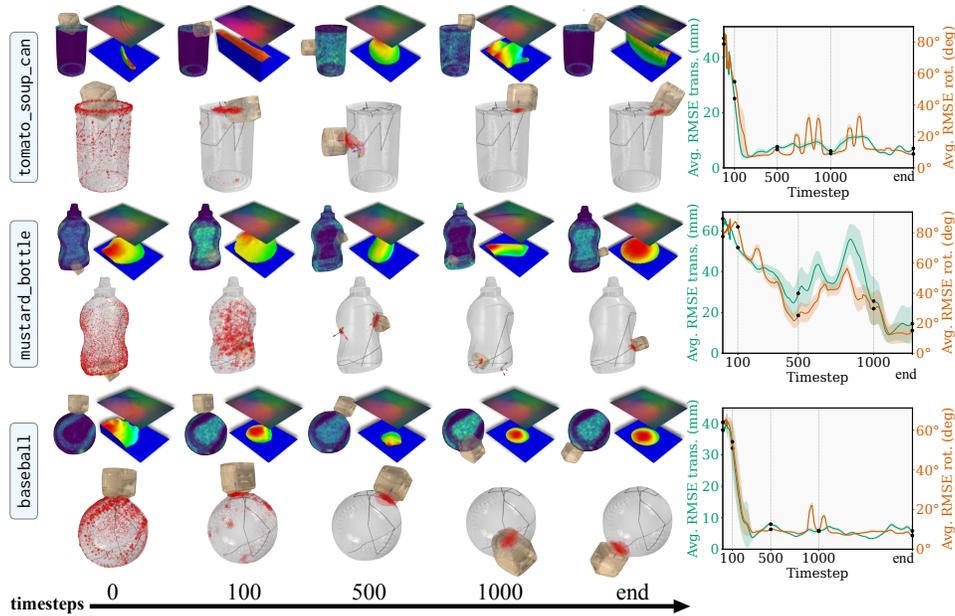


Figure 7: Snapshots of simulated sliding results on three YCB objects. For each row: **[top]** the tactile images, local geometries, and heatmap of pose likelihood with respect to the tactile codebook, **[bottom]** pose distribution evolving over time, and converging to the most-likely hypothesis after encountering salient geometries, **[right]** average translation and rotation RMSE of the distribution over time with variance over 10 trials.

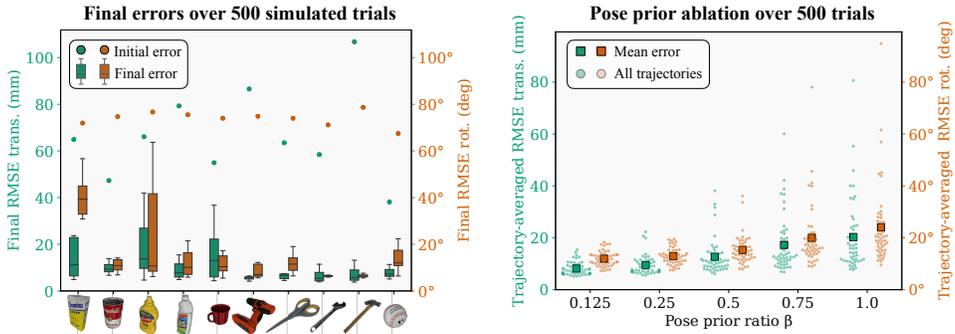


Figure 8: **[left]** Boxplot of final trajectory error over 500 simulation trials. For each object, we plot the averaged initial and final RMSE for the particle set. We observe better convergence for tools with salient geometries, as opposed to symmetric objects. **[right]** Ablation over initial pose uncertainty, to show lower average trajectory error with better visual priors. With a weaker initialization ($\beta = 1.0$), outliers in pose-error are more prevalent. tributed to tracking multiple modes that equally explain the same sliding sequence (please refer to supplementary video). We consider this a benefit of our multi-modal framework, and is especially prevalent for symmetric objects like the sugar_box, mustard_bottle, and mug. We perform ablations over pose prior (Figure 8 [right]), initializing each trial with uncertainty $\beta \times [\sigma_{\text{trans}}, \sigma_{\text{rot}}]$. This serves as a proxy for visual-priors; better initializations lead to fewer candidate modes. The full-stack operates, on average, at the rate of 10Hz.

Real-world: We run 500 similar trials for the real-world data from Section 5 with $\beta = 0.5$, which serves as a bound for a reasonable prior estimate available from vision. The real-world tactile heightmaps are noisier than simulation, thus a tighter initialization prevents particle depletion. In Figure 10 we present three representative trajectories that show the hypotheses converge to the ground-truth pose. The accumulated

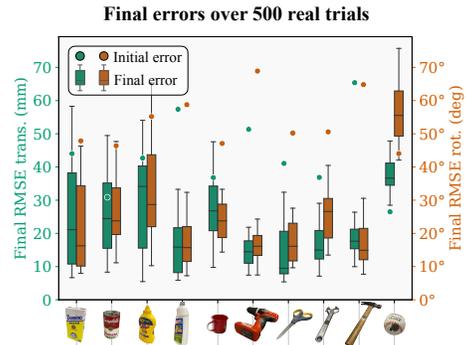


Figure 9: Boxplot of final error over 500 real-world trials from the 50 YCB-Slide trajectories.

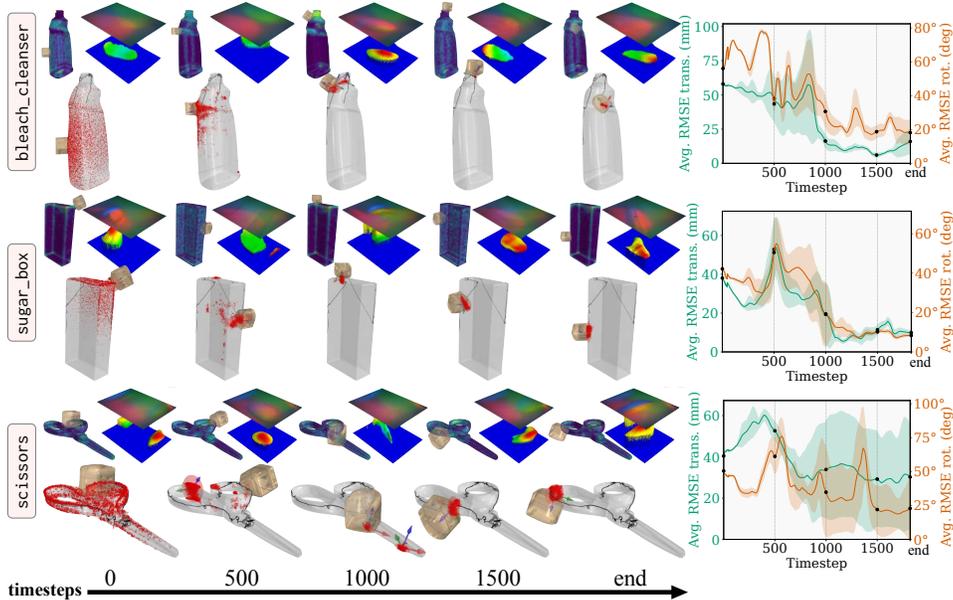


Figure 10: Snapshots of real-world sliding results on three YCB objects. For each row: **[top]** the tactile images, local geometries, and heatmap of pose likelihood with respect to the tactile codebook, **[bottom]** pose distribution evolving over time, and converging to the most-likely hypothesis after encountering salient geometries, **[right]** average translation and rotation RMSE of the distribution over time with variance over 10 trials.

statistics in Figure 9 shows reduced final error across all objects, except baseball. The averaged final pose errors are higher than the simulation results: 1.97cm and 21.48°. Once again, we see tools and intricate objects are the easiest to localize on while symmetric objects are the hardest. An anomaly is the baseball, on which we fail to localize in all trials. We attribute this to a lack of distinct geometry detected from real tactile images, effectively meaning we are trying to localize on a featureless sphere. These failure modes are highlighted in Figure 11. We present best-hypothesis error metrics and further qualitative results in Appendix E.

7 Conclusion and limitations

In this work, we demonstrate finger-object global localization from posed tactile images. This online method outputs a pose-distribution on the object’s surface that converges over time as the sensor traverses salient surface geometries. Specifically, our system is the first to learn tactile embeddings from local 3D geometry, and disambiguate them with a nonparametric particle filter. Our experiments demonstrate the surprising effectiveness of pairing learned tactile perception with Monte-Carlo methods to resolve distribution ambiguities.

Limitations: The current formulation is limited to a moving sensor relative to a fixed-pose object (or vice-versa). In future work, we wish to incorporate a dynamic object in our motion model through (i) visual measurements [26], and/or (ii) local in-hand tracking [20]. Our method currently does not work in scenarios where we lack ground-truth object models. Along with reconstructing objects [70], we would need to build the tactile codebook on-the-fly. For future in-hand manipulation tasks, it is also necessary to scale MidasTouch to multi-contact configurations. Our on-surface assumption leads to poor behavior when we break contact with the object, and the effects of shear on the tactile image sequence have been ignored [71]. Object compositionality doesn’t hold for deformable or articulated objects, so modeling these properties is an interesting future direction [72]. With a differentiable filter [73] we can fine-tune the system end-to-end for more robust real-world performance.

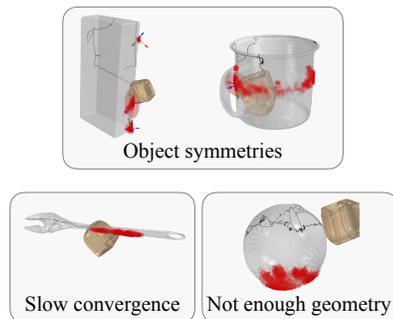


Figure 11: Failure modes in the real-world: (i) it may be hard to converge to the true hypothesis for objects with symmetries, (ii) slow convergence of the filter can lead to large pose uncertainty, (iii) lack of discernible geometry can result in drift from the true mode.

Acknowledgments

We thank Wenzhen Yuan, Ming-Fang Chang, Wei Dong, Maria Bauza Villalonga, and Antonia Bronars for insightful discussions. We are grateful towards the CMU AI Maker Space and Greg Armstrong for facilitating the collection of the YCB-Slide dataset. The authors acknowledge funding from Meta AI, and this work was partially carried out while Sudharshan Suresh interned at Meta AI.

References

- [1] J. Heller. *Catch-22: a novel*, volume 4. Simon and Schuster, 1961.
- [2] J. Shi, J. Z. Woodruff, P. B. Umbanhowar, and K. M. Lynch. Dynamic in-hand sliding manipulation. *IEEE Trans. on Robotics (TRO)*, 33(4):778–795, 2017.
- [3] N. Chavan-Dafle, R. Holladay, and A. Rodriguez. In-hand manipulation via motion cones. *Proc. Robotics: Science and Systems (RSS)*, 2018.
- [4] J. Shi, H. Weng, and K. M. Lynch. In-hand sliding regrasp with spring-sliding compliance and external constraints. *IEEE Access*, 8:88729–88744, 2020.
- [5] Y. She, S. Wang, S. Dong, N. Sunil, A. Rodriguez, and E. Adelson. Cable manipulation with a tactile-reactive gripper. *Intl. J. of Robotics Research (IJRR)*, 40(12-14):1385–1401, 2021.
- [6] J. Kerr, H. Huang, A. Wilcox, R. Hoque, J. Ichnowski, R. Calandra, and K. Goldberg. Learning self-supervised representations from vision and touch for active sliding perception of deformable surfaces. *arXiv preprint arXiv:2209.13042*, 2022.
- [7] O. M. Andrychowicz, B. Baker, M. Chociej, R. Jozefowicz, B. McGrew, J. Pachocki, A. Petron, M. Plappert, G. Powell, A. Ray, et al. Learning dexterous in-hand manipulation. *Intl. J. of Robotics Research (IJRR)*, 39(1):3–20, 2020.
- [8] T. Chen, J. Xu, and P. Agrawal. A system for general in-hand object re-orientation. In *Proc. Conf. on Robot Learning (CoRL)*, pages 297–307. PMLR, 2022.
- [9] A. Handa, A. Allshire, V. Makoviychuk, A. Petrenko, R. Singh, J. Liu, D. Makoviichuk, K. Van Wyk, A. Zhurkevich, B. Sundaralingam, Y. Narang, J.-F. Lafleche, D. Fox, and G. State. DeXtreme: Transfer of agile in-hand manipulation from simulation to reality. *arXiv*, 2022.
- [10] H. Qi, A. Kumar, R. Calandra, Y. Ma, and J. Malik. In-Hand Object Rotation via Rapid Motor Adaptation. In *Conference on Robot Learning (CoRL)*, 2022.
- [11] B. Sundaralingam and T. Hermans. Geometric in-hand regrasp planning: Alternating optimization of finger gaits and in-grasp manipulation. In *Proc. IEEE Intl. Conf. on Robotics and Automation (ICRA)*, pages 231–238. IEEE, 2018.
- [12] S. Brahmabhatt, A. Handa, J. Hays, and D. Fox. Contactgrasp: Functional multi-finger grasp synthesis from contact. In *Proc. IEEE/RSJ Intl. Conf. on Intelligent Robots and Systems (IROS)*, pages 2386–2393. IEEE, 2019.
- [13] W. Yuan, S. Dong, and E. H. Adelson. GelSight: High-resolution robot tactile sensors for estimating geometry and force. *Sensors*, 17(12):2762, 2017.
- [14] E. Donlon, S. Dong, M. Liu, J. Li, E. Adelson, and A. Rodriguez. GelSlim: A high-resolution, compact, robust, and calibrated tactile-sensing finger. In *Proc. IEEE/RSJ Intl. Conf. on Intelligent Robots and Systems (IROS)*, pages 1927–1934. IEEE, 2018.
- [15] M. Lambeta, P.-W. Chou, S. Tian, B. Yang, B. Maloon, V. R. Most, D. Stroud, R. Santos, A. Byagowi, G. Kammerer, et al. DIGIT: A novel design for a low-cost compact high-resolution tactile sensor with application to in-hand manipulation. *IEEE Robotics and Automation Letters (RA-L)*, 5(3):3838–3845, 2020.

- [16] 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 *Proc. IEEE/RSJ Intl. Conf. on Intelligent Robots and Systems (IROS)*, pages 3988–3993. IEEE, 2014.
- [17] M. Bauza, E. Valls, B. Lim, T. Sechopoulos, and A. Rodriguez. Tactile object pose estimation from the first touch with geometric contact rendering. In *Proc. Conf. on Robot Learning, CoRL*, 2020.
- [18] M. Bauza, A. Bronars, and A. Rodriguez. Tac2Pose: Tactile object pose estimation from the first touch. *arXiv preprint arXiv:2204.11701*, 2022.
- [19] P. Sodhi, M. Kaess, M. Mukadam, and S. Anderson. Learning tactile models for factor graph-based estimation. In *Proc. IEEE Intl. Conf. on Robotics and Automation (ICRA)*, pages 13686–13692. IEEE, 2021.
- [20] P. Sodhi, M. Kaess, M. Mukadam, and S. Anderson. Patchgraph: In-hand tactile tracking with learned surface normals. In *Proc. IEEE Intl. Conf. on Robotics and Automation (ICRA)*, 2022.
- [21] J. Hawkins, M. Lewis, M. Klukas, S. Purdy, and S. Ahmad. A framework for intelligence and cortical function based on grid cells in the neocortex. *Frontiers in neural circuits*, 12:121, 2019.
- [22] S. J. Lederman and R. L. Klatzky. Hand movements: A window into haptic object recognition. *Cognitive psychology*, 19(3):342–368, 1987.
- [23] F. Dellaert, D. Fox, W. Burgard, and S. Thrun. Monte Carlo localization for mobile robots. In *Proc. IEEE Intl. Conf. on Robotics and Automation (ICRA)*, volume 2, pages 1322–1328. IEEE, 1999.
- [24] B. Calli, A. Singh, J. Bruce, A. Walsman, K. Konolige, S. Srinivasa, P. Abbeel, and A. M. Dollar. Yale-CMU-Berkeley dataset for robotic manipulation research. *Intl. J. of Robotics Research (IJRR)*, 36(3):261–268, 2017.
- [25] L. Downs, A. Francis, N. Koenig, B. Kinman, R. Hickman, K. Reymann, T. B. McHugh, and V. Vanhoucke. Google scanned objects: A high-quality dataset of 3D scanned household items. *Proc. IEEE Intl. Conf. on Robotics and Automation (ICRA)*, 2022.
- [26] X. Deng, A. Mousavian, Y. Xiang, F. Xia, T. Bretl, and D. Fox. PoseRBPF: A Rao–Blackwellized particle filter for 6D object pose tracking. *IEEE Trans. on Robotics (TRO)*, 37(5):1328–1342, 2021.
- [27] A. Yamaguchi and C. G. Atkeson. Combining finger vision and optical tactile sensing: Reducing and handling errors while cutting vegetables. In *Proc. IEEE-RAS Intl. Conf. on Humanoid Robots (Humanoids)*, pages 1045–1051. IEEE, 2016.
- [28] B. Ward-Cherrier, N. Pestell, L. Cramphorn, B. Winstone, M. E. Giannaccini, J. Rossiter, and N. F. Lepora. The TacTip family: Soft optical tactile sensors with 3D-printed biomimetic morphologies. *Soft robotics*, 5(2):216–227, 2018.
- [29] A. Alspach, K. Hashimoto, N. Kuppawamy, and R. Tedrake. Soft-bubble: A highly compliant dense geometry tactile sensor for robot manipulation. In *Proc. IEEE Intl. Conf. on Soft Robotics (RoboSoft)*, pages 597–604. IEEE, 2019.
- [30] A. Padmanabha, F. Ebert, S. Tian, R. Calandra, C. Finn, and S. Levine. OmniTact: A multi-directional high-resolution touch sensor. In *Proc. IEEE Intl. Conf. on Robotics and Automation (ICRA)*, pages 618–624. IEEE, 2020.
- [31] S. Wang, Y. She, B. Romero, and E. Adelson. GelSight Wedge: Measuring high-resolution 3D contact geometry with a compact robot finger. In *Proc. IEEE Intl. Conf. on Robotics and Automation (ICRA)*. IEEE, 2021.

- [32] S. Wang, M. M. Lambeta, P.-W. Chou, and R. Calandra. TACTO: A fast, flexible, and open-source simulator for high-resolution vision-based tactile sensors. *IEEE Robotics and Automation Letters (RA-L)*, 2022.
- [33] A. Agarwal, T. Man, and W. Yuan. Simulation of vision-based tactile sensors using physics based rendering. In *Proc. IEEE Intl. Conf. on Robotics and Automation (ICRA)*, pages 1–7. IEEE, 2021.
- [34] Z. Si and W. Yuan. Taxim: An example-based simulation model for gelsight tactile sensors. *IEEE Robotics and Automation Letters (RA-L)*, 2022.
- [35] C. Choy, J. Park, and V. Koltun. Fully convolutional geometric features. In *Proc. Intl. Conf. on Computer Vision (ICCV)*, pages 8958–8966, 2019.
- [36] M. A. Uy and G. H. Lee. PointNetVLAD: Deep point cloud based retrieval for large-scale place recognition. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, pages 4470–4479, 2018.
- [37] J. Komorowski. MinkLoc3D: Point cloud based large-scale place recognition. In *Proc. Winter Conf. on Applications of Computer Vision (WACV)*, pages 1790–1799, 2021.
- [38] C. Corcoran and R. Platt. Tracking object pose and shape during robot manipulation based on tactile information. In *Proc. IEEE Intl. Conf. on Robotics and Automation (ICRA)*, volume 2. Citeseer, 2010.
- [39] A. Petrovskaya and O. Khatib. Global localization of objects via touch. *IEEE Trans. on Robotics (TRO)*, 27(3):569–585, 2011.
- [40] L. Zhang and J. C. Trinkle. The application of particle filtering to grasping acquisition with visual occlusion and tactile sensing. In *2012 IEEE International Conference on Robotics and Automation*, pages 3805–3812. IEEE, 2012.
- [41] M. Chalon, J. Reinecke, and M. Pfanne. Online in-hand object localization. In *Proc. IEEE/RSJ Intl. Conf. on Intelligent Robots and Systems (IROS)*, pages 2977–2984. IEEE, 2013.
- [42] B. Saund, S. Chen, and R. Simmons. Touch based localization of parts for high precision manufacturing. In *Proc. IEEE Intl. Conf. on Robotics and Automation (ICRA)*, pages 378–385. IEEE, 2017.
- [43] M. C. Koval, M. Klingensmith, S. S. Srinivasa, N. S. Pollard, and M. Kaess. The manifold particle filter for state estimation on high-dimensional implicit manifolds. In *Proc. IEEE Intl. Conf. on Robotics and Automation (ICRA)*, pages 4673–4680. IEEE, 2017.
- [44] S. Thrun. Particle filters in robotics. In *UAI*, volume 2, pages 511–518. Citeseer, 2002.
- [45] Z. Pezzementi, C. Reyda, and G. D. Hager. Object mapping, recognition, and localization from tactile geometry. In *2011 IEEE International Conference on Robotics and Automation*, pages 5942–5948. IEEE, 2011.
- [46] S. Luo, W. Mou, K. Althoefer, and H. Liu. Localizing the object contact through matching tactile features with visual map. In *Proc. IEEE Intl. Conf. on Robotics and Automation (ICRA)*, pages 3903–3908. IEEE, 2015.
- [47] J. Bimbo, S. Luo, K. Althoefer, and H. Liu. In-hand object pose estimation using covariance-based tactile to geometry matching. *IEEE Robotics and Automation Letters (RA-L)*, 1(1): 570–577, 2016.
- [48] K.-T. Yu and A. Rodriguez. Realtime state estimation with tactile and visual sensing. application to planar manipulation. In *Proc. IEEE Intl. Conf. on Robotics and Automation (ICRA)*, pages 7778–7785. IEEE, 2018.

- [49] A. N. Chaudhury, T. Man, W. Yuan, and C. G. Atkeson. Using collocated vision and tactile sensors for visual servoing and localization. *IEEE Robotics and Automation Letters (RA-L)*, 7(2):3427–3434, 2022.
- [50] S. Dong and A. Rodriguez. Tactile-based insertion for dense box-packing. In *Proc. IEEE/RSJ Intl. Conf. on Intelligent Robots and Systems (IROS)*, pages 7953–7960. IEEE, 2019.
- [51] S. Kim and A. Rodriguez. Active extrinsic contact sensing: Application to general peg-in-hole insertion. In *Proc. IEEE Intl. Conf. on Robotics and Automation (ICRA)*. IEEE, 2022.
- [52] T. Kelestemur, R. Platt, and T. Padir. Tactile pose estimation and policy learning for unknown object manipulation. *arXiv preprint arXiv:2203.10685*, 2022.
- [53] R. Gao, Z. Si, Y.-Y. Chang, S. Clarke, J. Bohg, L. Fei-Fei, W. Yuan, and J. Wu. ObjectFolder 2.0: A multisensory object dataset for Sim2Real transfer. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2022.
- [54] N. Kuppaswamy, 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 (RA-L)*, 5(2):1811–1818, 2019.
- [55] M. K. Johnson and E. H. Adelson. Retrographic sensing for the measurement of surface texture and shape. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, pages 1070–1077, 2009. doi:10.1109/CVPR.2009.5206534.
- [56] M. K. Johnson, F. Cole, A. Raj, and E. H. Adelson. Microgeometry capture using an elastomeric sensor. *ACM Trans. Graph.*, 30(4), July 2011. ISSN 0730-0301.
- [57] M. Bauza, O. Canal, and A. Rodriguez. Tactile mapping and localization from high-resolution tactile imprints. In *Proc. IEEE Intl. Conf. on Robotics and Automation (ICRA)*, pages 3811–3817. IEEE, 2019.
- [58] R. Ambrus, V. Guizilini, N. Kuppaswamy, A. Beaulieu, A. Gaidon, and A. Alspach. Monocular depth estimation for soft visuotactile sensors. In *Proc. IEEE Intl. Conf. on Soft Robotics (RoboSoft)*, 2021.
- [59] S. Suresh, Z. Si, J. G. Mangelson, W. Yuan, and M. Kaess. ShapeMap 3-D: Efficient shape mapping through dense touch and vision. In *Proc. IEEE Intl. Conf. on Robotics and Automation (ICRA)*, Philadelphia, PA, USA, May 2022.
- [60] C. R. Qi, H. Su, K. Mo, and L. J. Guibas. PointNet: Deep learning on point sets for 3D classification and segmentation. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, pages 652–660, 2017.
- [61] C. Choy, J. Gwak, and S. Savarese. 4D spatio-temporal ConvNets: Minkowski convolutional neural networks. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, pages 3075–3084, 2019.
- [62] I. Laina, C. Rupprecht, V. Belagiannis, F. Tombari, and N. Navab. Deeper depth prediction with fully convolutional residual networks. In *Proc. Intl. Conf. on 3D Vision (3DV)*, pages 239–248. IEEE, 2016.
- [63] T.-Y. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, and S. Belongie. Feature pyramid networks for object detection. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, pages 2117–2125, 2017.
- [64] L. Van der Maaten and G. Hinton. Visualizing data using t-sne. *Journal of machine learning research*, 9(11), 2008.

- [65] L. Pineda, T. Fan, M. Monge, S. Venkataraman, P. Sodhi, R. T. Chen, J. Ortiz, D. DeTone, A. Wang, S. Anderson, J. Dong, B. Amos, and M. Mukadam. Theseus: A Library for Differentiable Nonlinear Optimization. *Advances in Neural Information Processing Systems*, 2022.
- [66] W. Yuan, R. Li, M. A. Srinivasan, and E. H. Adelson. Measurement of shear and slip with a GelSight tactile sensor. In *Proc. IEEE Intl. Conf. on Robotics and Automation (ICRA)*, pages 304–311. IEEE, 2015.
- [67] S. Thrun. Probabilistic robotics. *Communications of the ACM*, 45(3):52–57, 2002.
- [68] E. Schubert, J. Sander, M. Ester, H. P. Kriegel, and X. Xu. DBSCAN revisited, revisited: why and how you should (still) use DBSCAN. *ACM Transactions on Database Systems (TODS)*, 42(3):1–21, 2017.
- [69] F. L. Markley, Y. Cheng, J. L. Crassidis, and Y. Oshman. Averaging quaternions. *Journal of Guidance, Control, and Dynamics*, 30(4):1193–1197, 2007.
- [70] S. Suresh, M. Bauza, K.-T. Yu, J. G. Mangelson, A. Rodriguez, and M. Kaess. Tactile SLAM: Real-time inference of shape and pose from planar pushing. In *Proc. IEEE Intl. Conf. on Robotics and Automation (ICRA)*, May 2021.
- [71] K. Aquilina, D. A. Barton, and N. F. Lepora. Shear-invariant sliding contact perception with a soft tactile sensor. In *Proc. IEEE Intl. Conf. on Robotics and Automation (ICRA)*, pages 4283–4289. IEEE, 2019.
- [72] Y. Wi, P. Florence, A. Zeng, and N. Fazeli. VIRDO: Visio-tactile implicit representations of deformable objects. In *Proc. IEEE Intl. Conf. on Robotics and Automation (ICRA)*. IEEE, 2022.
- [73] R. Jonschkowski, D. Rastogi, and O. Brock. Differentiable particle filters: End-to-end learning with algorithmic priors. In *Proc. Robotics: Science and Systems (RSS)*, 2018.
- [74] McMaster-Carr. <https://www.mcmaster.com/>. Accessed: 2022-08-25.

Appendices

Appendix A Tactile depth network: data and training

The TDN is trained with TACTO [32] images, minimizing heightmap reconstruction loss, as in the monocular depth estimation [62]. These image-heightmap pairs are shown in Figure 12 [right]. The TACTO interactions are over a diverse set of YCB objects shown in Figure 12 [left].

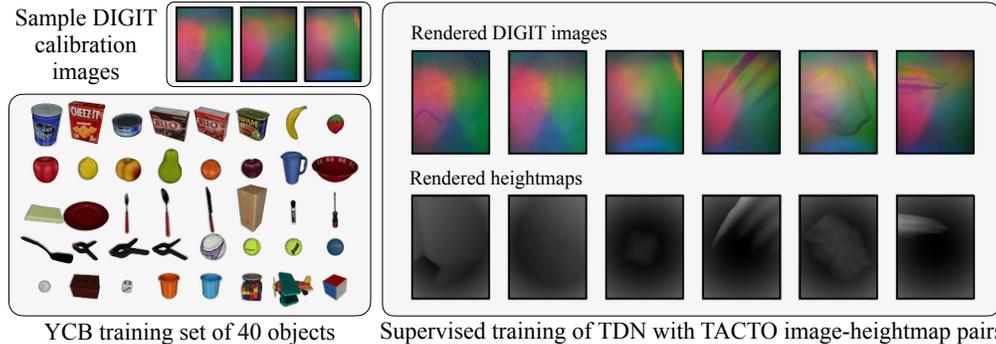


Figure 12: [top-left] Three calibration images captured from different DIGIT sensors. We augment our training data with these calibrations for better sim2real transfer. [bottom-left] The 40 YCB [24] objects across food, kitchen, tool, shape, and task classes. [right] Examples of image-heightmap pairs used for supervision.

Generating contact poses: We render realistic DIGIT images from 5000 unique poses on each object. To generate these contact poses, we sample point-normal pairs across the object’s mesh. Through rejection-sampling, we can get an approximately even distribution across the surface. Additionally, we prioritize sampling edges with mesh feature angles $> 10^\circ$. This gives us the uniform spread that we desire, while also capturing the salient features across the object classes.

After sampling these points, we add penetration depth randomly sampled between 0.5mm to 2mm. To convert a contact location to a pose, we first assign a random orientation angle ϕ around the surface normal direction between $[0^\circ, 360^\circ]$. We add an orientation noise angle $\theta = \mathcal{N}(0, 5^\circ)$ in the cone perpendicular to the surface normal to ensure that the poses aren’t always orthogonal to the local surface. We randomly assign 2% of all poses to not make contact with the surface.

Image augmentations: For sim2real transfer, the TDN should generalize across different sensor lighting conditions. For example, in our real experiments we use three different DIGIT sensors, and there is some wear and tear of the elastomer over time. Through TACTO, we can calibrate the renderer with respect to real-world images. These images, pictured in Figure 12 [left], are captured when the DIGIT does not make contact with a surface. We collect 10 calibration images over the YCB-Slide dataset, and use each as the calibration for a subset of the training data generation.

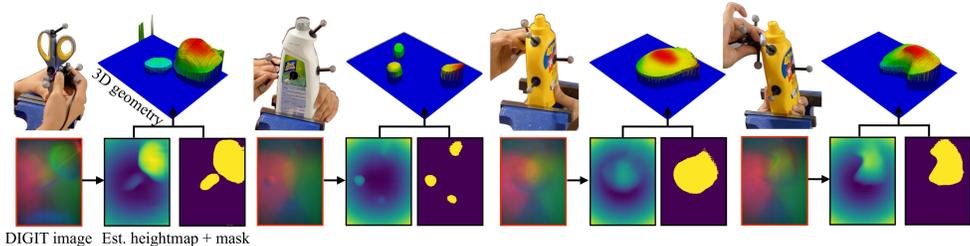


Figure 13: Extended results from the tactile depth network, similar to Figure 2. Given input image, the network predicts heightmap which is reprojected to give 3D geometry.

Appendix B Tactile codes: 3D versus image

In this section, we compare our 3D tactile codes against a baseline tactile image embedding. While Bauza et al. [18] use a similar contrastive strategy to learn image embeddings, they are object-specific. We choose instead to compare against the learned model in ObjectFolder 2.0 [53]. They extract features from the fully-convolutional residual network bottleneck layer, the tactile depth

network we use as our observation model (Section 4.1). They use these embeddings for multi-touch contact location estimation with GelSight images.

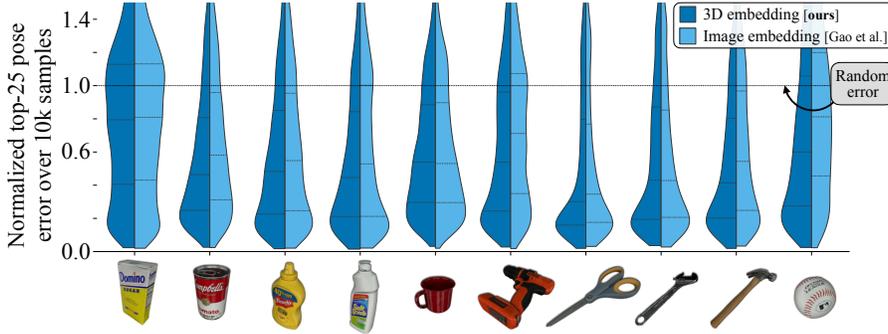


Figure 14: Pose-error for 10k single-touch queries on YCB objects, comparing our 3D tactile codes v.s. image embeddings. For each query, we get the top-25 highest scores from the tactile codebook C_o , and compute their minimum pose-error with respect to ground-truth. This density distribution is plotted as a violin-plot, normalized by the error from a randomly-sampled touch.

We perform the same single-touch localization experiments from Figure 4, using both our tactile codes, and the image embeddings. For each method, we build an object-specific codebook. For each query tactile image, we generate its corresponding embedding and match against the codebook. We then compute the minimum pose-error from the top-25 matches, and repeat this for 10k touch queries. From Figure 14 we observe lower pose-errors for our tactile codes, with an average normalized pose-error of **0.473** compared to 0.546 from Gao et al. [53]. Also importantly, our embeddings are low-dimensional, leading to a codebook size 300 times smaller than that of the baseline. Examples of top single-touch query results for sugar_box are shown in Figure 15.

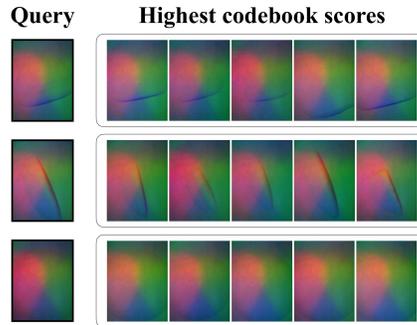


Figure 15: Examples of query tactile images matched against the tactile codebook of the sugar_box object. We see that the top scores in each case are images whose sensor pose is also similar to the query.

Appendix C Particle filter: Implementation details

Particle count: We initialize liberally, with $N_0 = 50k$ so as to better capture poses close to the ground-truth. With too few particles, we run the risk of not capturing good 6D pose candidates initially. In practice, reducing N_t greatly improves computation time, but is a trade-off on tracking accuracy. Alongside resampling, we dynamically adjust the number of particles N_t based on filter convergence. We track the average standard deviation of the hypothesis set σ_{h_t} , and inject or remove particles according to the ratio $\sigma_{h_t}/\sigma_{h_{t-1}}$. When injecting particles, we replicate those with the largest weights, and when removing particles, we delete those with the lowest. To prevent particle depletion, we do not let it drop below 1k particles.

Pruning: We leverage our on-surface assumption, to prune particles that drift too far away from the objects surface. Given the object meshes, we construct a k-d tree of all vertices, and at each iteration we perform a nearest neighbor search for the particles using this tree. When the distance check exceeds 2mm, which we consider the sensor’s maximum penetration distance, we set the corresponding particle’s weight to zero. This "soft" on-surface constraint allows for some noise in the odometry, while gradually pruning candidate hypotheses.

Real-world experiment: With real DIGIT images, we’ve found exponential time smoothing over predicted heightmaps gives stable local geometry and removes outlier effects. Further, we reduce the frequency of particle resampling to every five iterations. This prevents particle depletion in the presence of erroneous heightmaps. In the real-world experiments, there are instances where the DIGIT can slide off the object’s surface. In those cases, the TDN does not produce a point-cloud and we instead generate a randomized \mathbb{R}^{256} embedding.

Appendix D Additional YCB-Slide details



Figure 16: [left] Our data-collection setup with the YCB object, motion capture, DIGIT sensor, and recording camera. [right] The 10 YCB test objects during interactions, with representative tactile images.

Our real-world dataset was collected in the indoor environment pictured in Figure 16. Each object is tightly clamped onto a heavy-duty bench vise, while the DIGIT sensor is slid across its surface. The full set of 50 simulated interactions are shown in Figure 17. Each trajectory has a fixed geodesic length of approximately 0.5m.

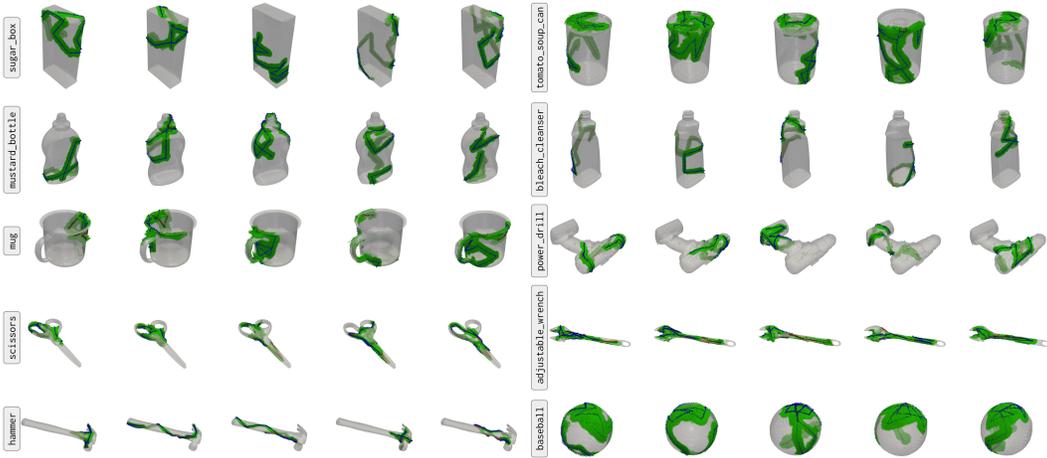


Figure 17: All 50 sliding trajectories from tactile simulation on the 10 YCB objects. Overlaid in green are the local 3D geometries captured by the tactile sensor, and the contact poses as RGB coordinate axes.

The 50 collected real-world interactions are shown in Figure 18. While these human-designed sequences are not random, each covers different sections of the object geometry.

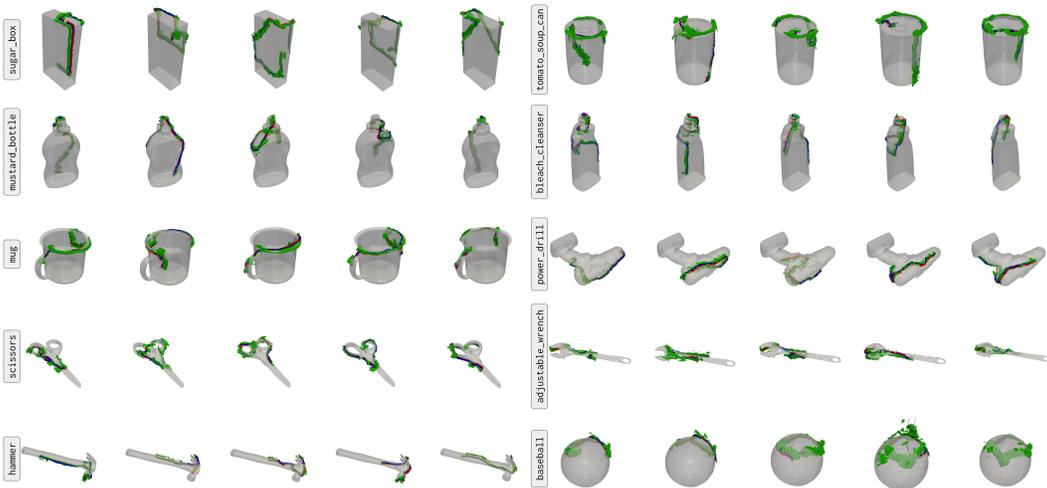


Figure 18: All 50 sliding trajectories from real-world interactions on the 10 YCB objects. Overlaid in green are the local 3D geometries captured by the tactile sensor, and the contact poses as RGB coordinate axes.

Appendix E Additional MidasTouch results

Closest hypothesis error: In Section 6, we present final RMSE for the pose particles with respect to ground truth. In Figure 8 and 9 we plot these accumulative statistics at the final timestep T :

$$e_{\text{trans}} = \sqrt{\frac{1}{N_T} \sum_{\mathbf{x} \in \mathbf{X}_T} \|\mathbf{x}_{\text{trans}} - \mathbf{x}_{\text{trans}}^{\text{gt}}\|_2^2}, \quad e_{\text{rot}} = \sqrt{\frac{1}{N_T} \sum_{\mathbf{x} \in \mathbf{X}_T} \|\mathbf{x}_{\text{rot}} - \mathbf{x}_{\text{rot}}^{\text{gt}}\|_2^2} \quad (3)$$

This is a general error metric for the particle filter, but penalizes a multi-modal pose distribution. We present an additional metric that computes RMSE with respect to hypothesis set h_T , and uses the error associated with the closest hypothesis to ground-truth:

$$\text{min_cluster}(e_{\text{trans}}) = \sqrt{\min_{\mathbf{x} \in h_T} \|\mathbf{x}_{\text{trans}} - \mathbf{x}_{\text{trans}}^{\text{gt}}\|_2^2}, \quad \text{min_cluster}(e_{\text{rot}}) = \sqrt{\min_{\mathbf{x} \in h_T} \|\mathbf{x}_{\text{rot}} - \mathbf{x}_{\text{rot}}^{\text{gt}}\|_2^2} \quad (4)$$

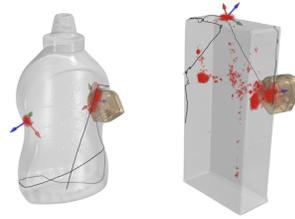


Figure 19: Examples of cases where particle pose error is high, but we still capture the true pose in our multi-modal distribution.

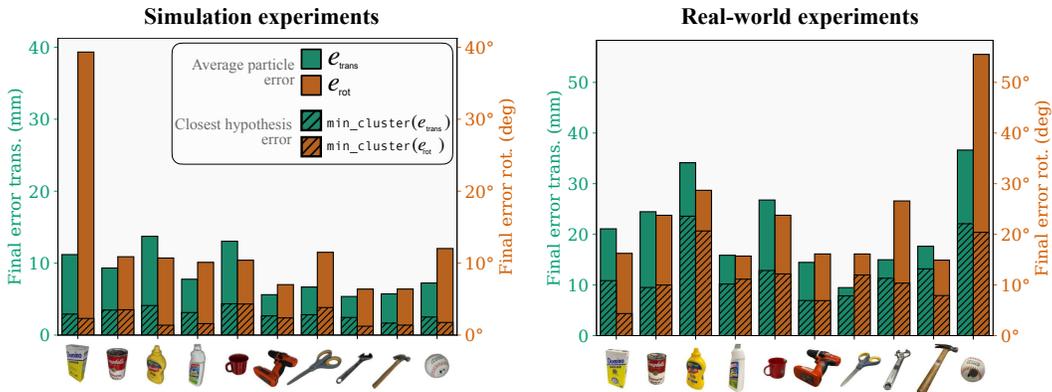


Figure 20: Plots for the 500 simulated [left] and real-world [right] trials comparing e_{trans} , e_{rot} against $\text{min_cluster}(e_{\text{trans}})$, $\text{min_cluster}(e_{\text{rot}})$. These serve as a complement to Figure 8 and 9, and highlight the multi-modality of the filtering problem. Given knowledge of ground-truth, we pick the cluster closest to it and plot the RMSE statistics with respect to it (Equation 4). We observe lower errors across all objects in both simulated and real settings, empirically indicating the true mode is captured in most cases.

While this assumes we have access to the ground-truth, it can better demonstrate if we capture the true pose in our multi-modal distribution. In Figure 20, we see the min_cluster statistics plotted alongside the original final RMSE. Across all experiments, we end up with lower error: with a median of $0.28\text{cm} + 2.03^\circ$ in simulation and $1.11\text{cm} + 10.76^\circ$ in the real-world.

Further qualitative results: Finally, we highlight some visualization similar to Figure 7 and 10, for the remaining YCB test objects. In Figures 21 and 22 we show snapshots of MidasTouch on the remaining YCB objects. Alongside the snapshots is the translation and rotation RMSE over time, averaged over 10 trials. We see filter convergence across different YCB objects, along with the failure mode of the baseball in Figure 22. Please refer to the supplementary video for further visualizations.

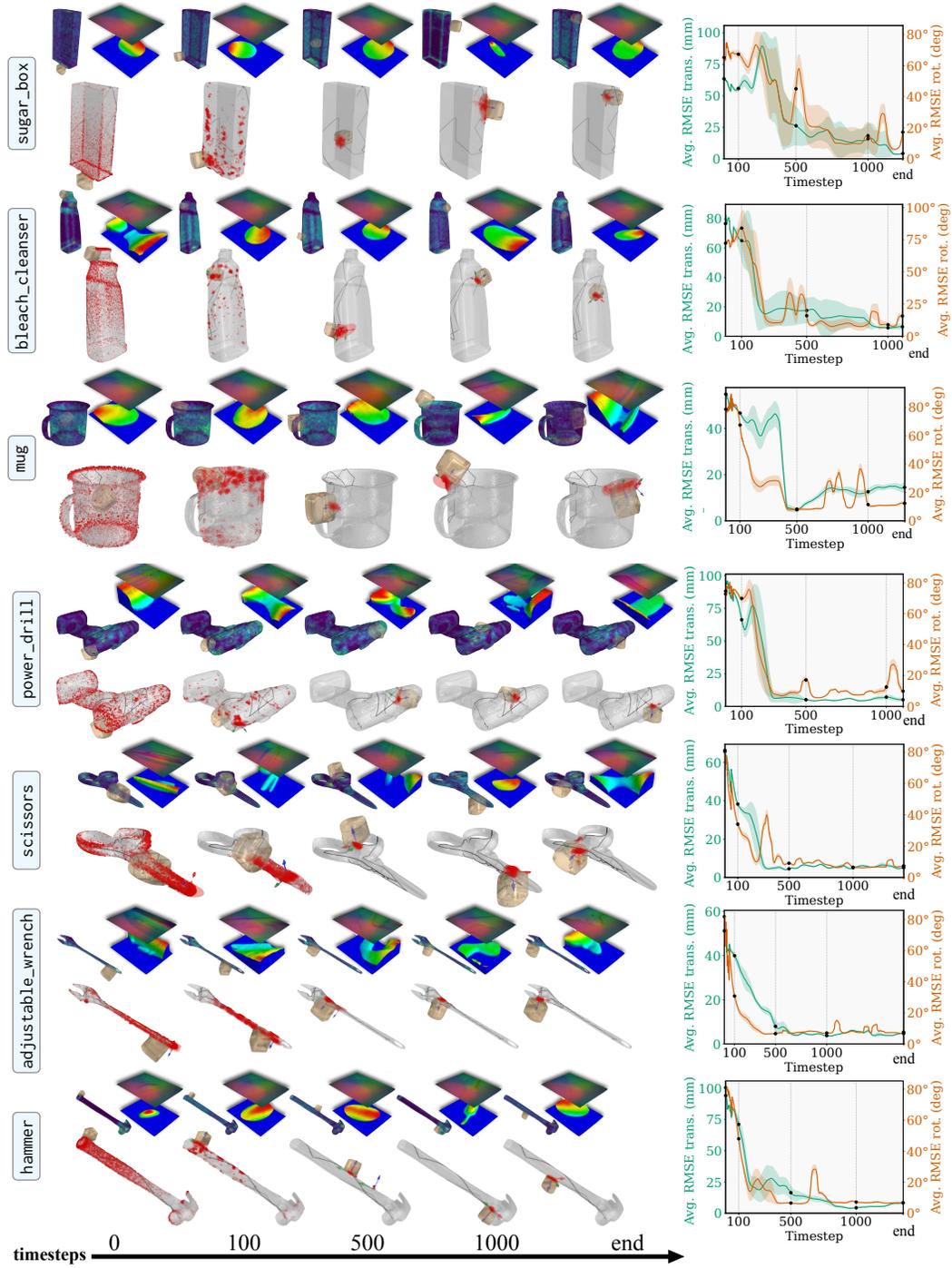


Figure 21: Select simulation results from seven objects not shown in Figure 7. For each row: **[top]** the tactile images, local geometries, and heatmap of pose likelihood with respect to the tactile codebook, **[bottom]** pose distribution evolving over time, **[right]** average translation/rotation RMSE of the distribution over time with variance over 10 trials.

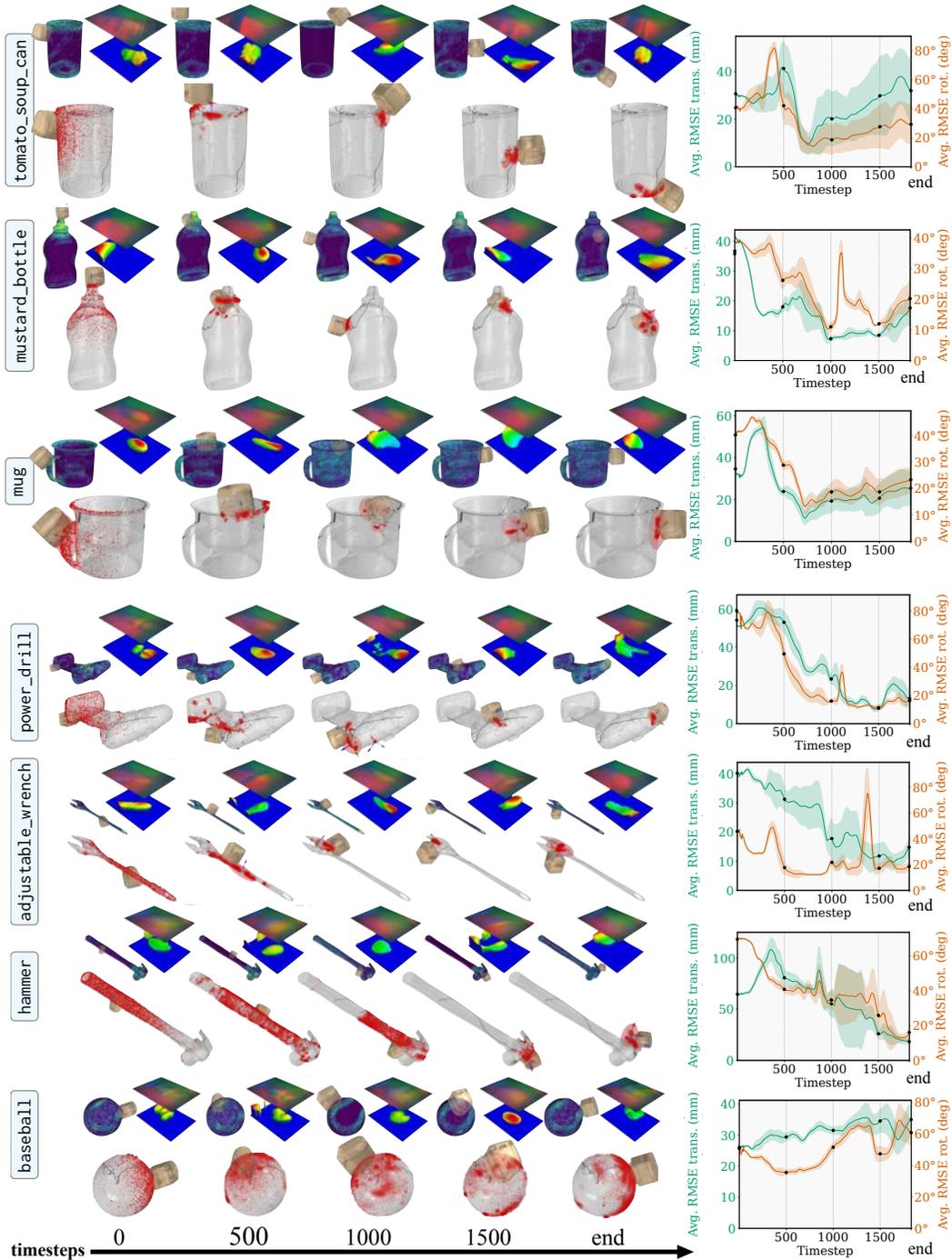


Figure 22: Select real-world results from seven objects not shown in Figure 10. For each row: **[top]** the tactile images, local geometries, and heatmap of pose likelihood with respect to the tactile codebook, **[bottom]** pose distribution evolving over time, **[right]** average translation/rotation RMSE of the distribution over time with variance over 10 trials.

Appendix F Study on contact patch area

We analyze the correlation of surface contact patch area with the performance of our filter. During interaction, it is crucial to maintain forceful contact with the surface area impinging the sensor. This gives us a larger contact area, and more 3-D surface geometry to match against the tactile codebook. For the DIGIT, this is theoretically between 0 to 6 cm², and can be obtained as the pixel area of C_t (Section 4.1).

To show the importance of larger contact areas, we record a single simulated trajectory on `power_drill`, ablated over five different penetration depth ranges. We randomly sample penetration depth in the range of $\omega \times \mathcal{N}(0.5, 2 \text{ mm})$, where $\omega \in [0.1, 0.325, 0.55, 0.775, 1.0]$. Figure 24 shows the same interaction with five different ω values. We observe that large penetration captures more surface geometry.

For each profile, we average the results of 10 filtering trials, and plot the final pose error v.s. average trajectory contact area. Figure 23 shows that more 3-D surface geometry can lead to lower downstream error in finger pose tracking. Intuitively, this is analogous to a depth-camera with a larger depth range, generating more complete scans of the scenes it perceives.

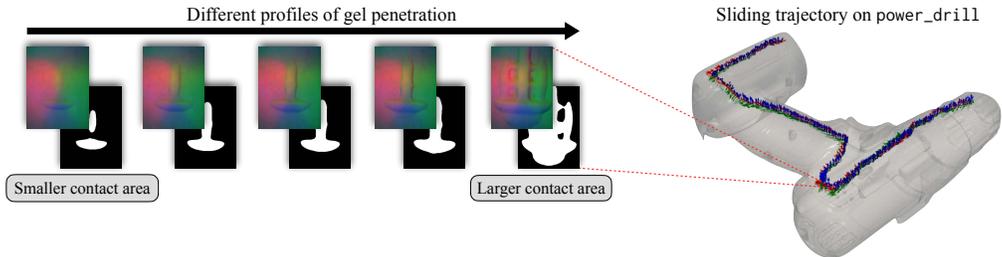


Figure 24: [right] Fixed trajectory on `power_drill` for which we apply different penetration profiles. [left] For the same local surface, different penetration profiles observe very different contact shapes and tactile images.

Appendix G Experiments on small parts

MidasTouch focuses on YCB-sized objects that we encounter in household and assistive robotics contexts. These are considerably larger than the robot finger, which is relevant for our desired applications of in-hand and tabletop manipulation. YCB-Slide spans objects with surface areas ranging from 109cm² (`adjustable_wrench`) to 643cm² (`bleach_cleanser`), while the sensor has a footprint of 6cm². In this section, we show simulated experiments for small parts with large sensor-model overlap, similar to prior work [16, 18].

We select three objects from McMaster-Carr [74], the `cotter_pin`, `eyebolt`, and `steel_nail`, each of 2" length. For each we generate a tactile codebook, and record a short simulated trajectory along the object's length (just as in Section 5). In Figure 26 we show results for all three, where the filter quickly converges to the true mode. We run each experiment 10 times, and show the accumulative statistics in Figure 25. We observe a final error of $\approx [4 \text{ mm}, 5^\circ]$, which is roughly twice as accurate as results in Section 6. Moreover, this requires trajectories 10 \times smaller, with 5 \times less particles. This is due to the small size of objects, and larger relative field-of-view.

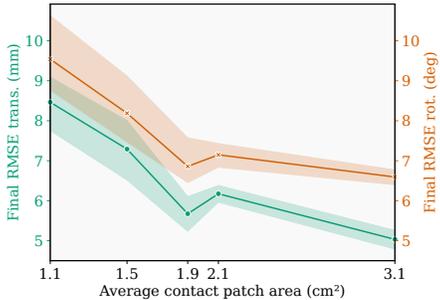


Figure 23: Plot over 50 trials of the pose error v.s. average contact area of the trajectory; larger contact areas lead to lower downstream error in tracking.

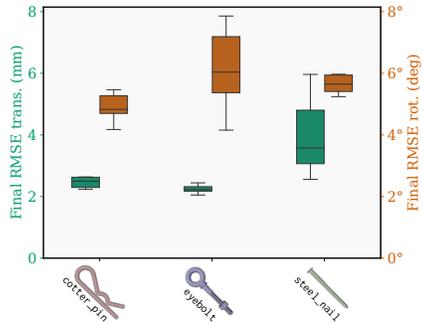


Figure 25: Boxplot of final error over 30 simulated trials on the McMaster-Carr small parts.

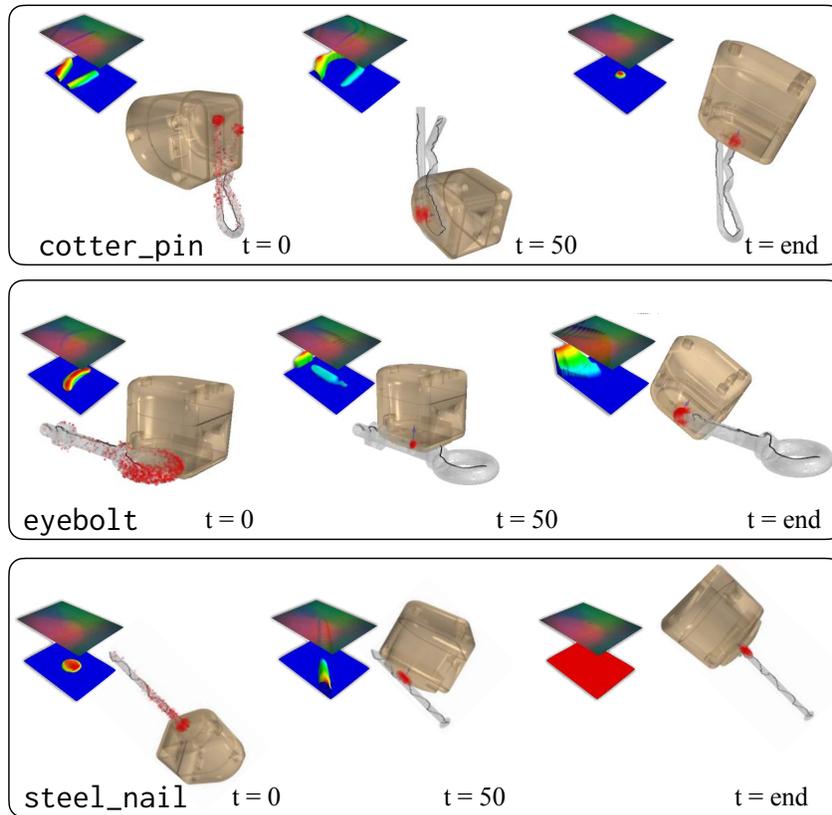


Figure 26: Filtering results from the three McMaster-Carr small parts. We visualize the tactile images, local geometries, and pose distribution evolving over time.