

Active Visuo-Haptic Object Shape Completion

Lukas Rustler¹, Jens Lundell², Jan Kristof Behrens³, Ville Kyrki², Matej Hoffmann¹

Abstract—Recent advancements in object shape completion have enabled impressive object reconstructions using only visual input. However, due to self-occlusion, the reconstructions have high uncertainty in the occluded object parts, which negatively impacts the performance of downstream robotic tasks such as grasping. In this work, we propose an active visuo-haptic shape completion method called Act-VH that actively computes where to touch the objects based on the reconstruction uncertainty. Act-VH reconstructs objects from point clouds and calculates the reconstruction uncertainty using IGR, a recent state-of-the-art implicit surface deep neural network. We experimentally evaluate the reconstruction accuracy of Act-VH against five baselines in simulation and in the real world. We also propose a new simulation environment for this purpose. The results show that Act-VH outperforms all baselines and that an uncertainty-driven haptic exploration policy leads to higher reconstruction accuracy than a random policy and a policy driven by Gaussian Process Implicit Surfaces. As a final experiment, we evaluate Act-VH and the best reconstruction baseline on grasping 10 novel objects. The results show that Act-VH reaches a significantly higher grasp success rate than the baseline on all objects. Together, this work opens up the door for using active visuo-haptic shape completion in more complex cluttered scenes.

Index Terms—Perception for Grasping and Manipulation; RGB-D Perception; Deep Learning for Visual Perception.

I. INTRODUCTION

SHAPe completion, that is reconstructing the shape of an object based on incomplete sensory information, is an active research problem with many potential applications in medicine and robotics. To date, most methods have reconstructed objects from only visual data, including RGB images, depth images, or point clouds. The main drawback of visual data is that it is incomplete as the objects self-occlude, *i.e.*, only the front side is visible from a single viewpoint.

Manuscript received: September, 9, 2021; Revised November, 19, 2021; Accepted January, 31, 2022.

This paper was recommended for publication by Editor Markus Vincze upon evaluation of the Associate Editor and Reviewers' comments. This work was supported by the project Interactive Perception-Action-Learning for Modelling Objects (IPALM) (H2020 – FET – ERA-NET Cofund – CHIST-ERA III / Technology Agency of the Czech Republic, EPSILON, no. TH05020001 / Academy of Finland, no. 326304). M.H. and L.R. were additionally supported by OP VVV MEYS funded project CZ.02.1.01/0.0/0.0/16_019/0000765 “Research Center for Informatics”. L.R. was also supported by the Czech Technical University in Prague, grant no. SGS20/128/OHK3/2T/13. J.K.B. was supported by the European Regional Development Fund under project Robotics for Industry 4.0 (reg. no. CZ.02.1.01/0.0/0.0/15_003/0000470).

¹ Lukas Rustler and Matej Hoffmann are with the Department of Cybernetics, Faculty of Electrical Engineering, CTU in Prague matej.hoffmann@fel.cvut.cz

² Jens Lundell and Ville Kyrki are with the Intelligent Robotics Group, Department of Electrical Engineering and Automation, School of Electrical Engineering, Aalto University, 02150 Espoo, Finland

³Jan Kristof Behrens is with the Czech Institute of Informatics, Robotics, and Cybernetics, CTU in Prague

Digital Object Identifier (DOI): see top of this page.

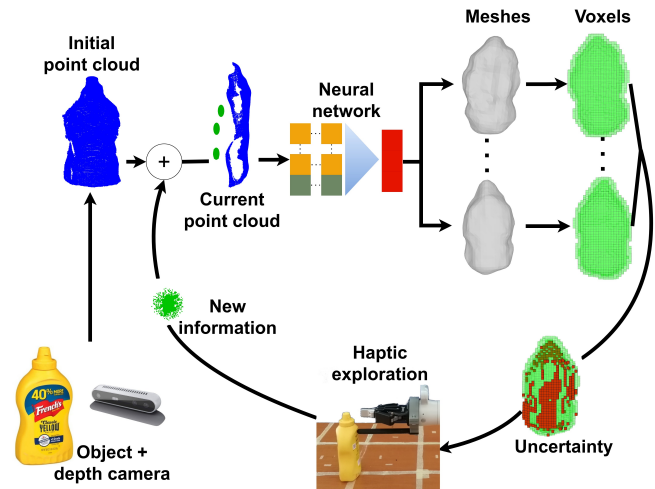


Fig. 1: Schematic operation of Act-VH. See text for details.

This increases the reconstruction uncertainty, which can negatively affect downstream robotic tasks such as grasping. A straightforward approach to combat the perceptual uncertainty is to gather additional data of the unseen object parts by touching the object and then reconstruct the object shape from combined visuo-haptic data. However, current visuo-haptic shape completion methods use heuristics to choose where to explore the objects [1] or require an impractical number of touches for a good reconstruction [2], [3].

We combat these issues in this work with Act-VH, a data-efficient closed-loop active visuo-haptic shape completion method. The operation is schematically illustrated in Fig. 1. First, a depth image is used to construct an initial point cloud. Then, a deep implicit surface network (IGR [4]) generates several possible shape reconstructions by iteratively refining randomly initialized latent codes until the reconstructed shapes fit the input point cloud. The discrepancy between the reconstructed shapes are then used to form a single voxel-grid reconstruction with uncertainty. The voxel with the highest uncertainty is selected for haptic exploration, adding a new point to the object representation. This process is repeated, further refining the shape reconstruction.

We experimentally validated Act-VH in simulation and the real world. To validate the method in simulation, we developed a new visuo-haptic benchmark task. Using this task, we compared the reconstruction accuracy of Act-VH with random and uncertainty-driven exploration to five baselines on 105 reconstructions. The results showed that Act-VH with uncertainty-driven exploration outperformed the baselines. In real-world experiments, we validated Act-VH in terms of reconstruction accuracy and grasp success rates on 10 objects. Similar to the simulation results, the real-world reconstruction

results also showed a significantly better accuracy using Act-VH. Finally, the grasping experiment on the real robot showed that the grasp success rates after 5 touches increased from 30% to 80% using Act-VH, while for the best reconstruction baseline it only increased from 20% to 46.7%, once again showing the benefits of Act-VH.

The main contributions of this work are: (i) a novel active visuo-haptic shape completion method; (ii) a visuo-haptic simulation environment; and (iii) an empirical evaluation of the proposed method against the state of the art, presenting improvements in terms of reconstruction accuracy (in simulation and on the real setup) and grasp success rates (real robot). The simulation environment, which at the same time serves as a benchmark, and the data from experiments are available at <https://github.com/ctu-vras/visuo-haptic-shape-completion>. An accompanying video is here: <https://youtu.be/iZF4ph4zMEA>.

II. RELATED WORK

An object shape can be reconstructed from visual input, haptic input, or their combination. Therefore, in this section, we split the review based on the sensory input used for reconstruction. Furthermore, the inputs can be collected only once or gathered actively to improve the reconstruction—the latter typically known as active perception [5]. Thus, we also review active reconstruction methods per input modality.

A. Visual-Only Shape Completion

Completing object shapes from visual data is the most common approach because the data capture global information about the object. Early visual shape completion approaches were geometry- or template-based. Examples of geometry-based approaches reconstruct objects by mirroring them through their symmetry axis [6] or using heuristics to fit primitives to resemble the object [7]. Template-based approaches search in a database for an object most similar to the perceived one [8]. The limitation of both methods is that they do not generalize well beyond specific objects. For instance, mirroring-based approaches result in poor reconstruction if the object has more than one axis of symmetry, while template matching will fail if the match is incorrect or no similar object exists in the database.

To combat the limitations of geometry- and template-based methods, Machine Learning (ML)-based shape completion approaches were proposed [1], [4], [9]–[14]. An early such approach trained a Gaussian Process Implicit Surface (GPIS) to reconstruct objects [9]. However, the GPIS reconstructed overly smooth objects and, due to its poor scaling to many data points, the input point cloud had to be down-sampled, losing valuable information. More recent ML approaches use Deep Learning (DL) techniques to train 3D Convolutional Neural Networks (CNNs) to complete the shape of objects represented as voxel grids [1], [10]–[12]. The limitation of voxel-based approaches is that the computation and memory requirements grow cubically with the object shape resolution. As such, fine object details are not preserved, which is essential when, for instance, sampling grasp proposals.

To overcome the issue with voxel-grid representations, researchers proposed new network architectures that can handle continuous shapes [4], [13], [14]. The architectures based on implicit surfaces are more computation- and memory-efficient than voxel-based representations and produce higher quality reconstructions. Because of these benefits, we chose to reconstruct objects with the implicit surface method IGR [4].

Despite the impressive results of visual-only shape completion, the noise in the visual data and the objects' self-occlusions result in high reconstruction uncertainty, especially on the nonvisible parts of the object. If there is a possibility to move the camera, these limitations can be alleviated by actively choosing alternative viewpoints (also called next-best-view) [15]. However, if the camera is not movable, another option is to use haptic data gathered by a robot.

B. Haptic-Only Shape Completion

If the robot has means to accurately detect and localize contacts with the object, tactile exploration can be more precise than visual data. Furthermore, any reachable part of the object—like its back side—can be explored. Most recent haptic-only shape completion approaches mainly reconstruct objects using classical ML models such as implicit shape potentials [16], Gaussian Processes (GPs) [17], GPIS's [18] or Gaussian Process Implicit Shape Potentials (GPISPs) [19]. Additionally, some haptic exploration approaches actively explore the object to reduce the uncertainty in the reconstruction [17]–[19]. The limitation with haptic data is its local nature—one touch only explores a small object region. Consequently, accurate object reconstruction from tactile data requires tens [19] to hundreds [17] of touches which is impractical for real robotic systems.

C. Visuo-Haptic Shape Completion

To address the limitations of visual- or haptic-only shape completion, some works have proposed visuo-haptic shape completion [2], [3], [20]–[24]. Most of these works reconstruct objects using ML techniques such as GPIS's [20], [21], GPs [3], CNNs [2], [22], or Graph Neural Networks (GNNs) [23], [24]. A limitation of the non-DL-based visuo-haptic approaches [3], [20], [21] is that good reconstructions often require haptic data all around the object. On the other hand, the main limitation of DL-based CNN approaches [2], [22] is the low object resolution, while for GNNs [23], [24] it is the non-smooth shape reconstruction and that the reconstructions are only evaluated in simulation.

Another known problem for all visuo-haptic shape completion works is deciding where to explore the object haptically. One solution is to use heuristics, such as always approaching the object directly opposite the camera [22]; another is to explore randomly [23]. Neither of these are particularly efficient as there exist more information-rich places to explore the object. To this end, some approaches learn where to explore the object [24] or use uncertainty of the reconstructions to guide exploration [2], [3], [20], [21].

The work presented here also does uncertainty-driven visuo-haptic shape completion using DL. Compared to similar works

that use CNNs [2], [22], we use implicit surface networks to reconstruct highly detailed and smooth objects. Compared to works using GNNs [23], [24], we evaluate our approach not only in simulation but also on real world reconstruction tasks. Furthermore, we propose a novel DL-based uncertainty-driven exploration strategy and evaluate if our method benefits robotic grasping.

III. METHOD

We propose the method in Fig. 1 to do active visuo-haptic shape completion. We assume that the visual measurements are only captured once while the haptic measurements are collected incrementally by exploring the object. It is assumed that the object does not move after haptic exploration. Based on these assumptions, the objective is to select a sequence of touches that would lead to the greatest improvements in reconstruction accuracy.

A. Uncertainty-Driven Haptic Exploration

Completing the shape of an object O perfectly from real world measurements Y is impossible due to the inherent noise and incompleteness of such measurements. The object O can be modeled in several ways. In this work, we specifically use multiple of these representations as shown in Fig. 1— from input point cloud, to Signed Distance Function (SDF) as used by the IGR, to mesh, and finally voxel grid, where the uncertainty is computed.

We propose to model the object O probabilistically as

$$P(O|Y), \quad (1)$$

where O represents the occupancy of the object and Y represents sensor measurements. The occupancy is represented as a voxel grid $O = (O^k)$ where k is the index of a voxel such that $P(O^k)$ is the probability that voxel k is part of the object. In this work, Y consists of visual v and haptic h data.

Formally, the objective is to, at each time step t , choose a location for haptic exploration that minimizes the uncertainty about the occupancy quantified as its variance

$$\operatorname{argmin}_{h_t \in H} \operatorname{Var}(O_t|v, h_{1:t-1}, h_t), \quad (2)$$

where $h_{1:t-1}$ is the data from previously executed haptic explorations and H is the set of all possible haptic explorations. Note that the variance at time 0 is based on visual data only $\operatorname{Var}(O_0|v)$.

Minimizing Eq. 2 requires a probabilistic model of the object's 3D shape, which is complex to form due to the high-dimensional nature of the data. Instead, we choose to approximate the model with a set of shape samples $o^{1:S}$ drawn from an underlying generative shape distribution $P(O_t|v, h_{1:t-1})$. The actual sampling process $o^s \sim P(O_t|v, h_{1:t-1})$ is described in the next section.

Assuming a set of shape samples are given in the form of voxel grids, we define the haptic exploration h_t that minimizes Eq. 2 to be the voxel k with the largest variance. This is formally expressed as

$$\operatorname{argmax}_{k \in K} \operatorname{Var}(O^k), \quad (3)$$

where k is a single voxel in a voxel grid K , and $\operatorname{Var}(O^k)$ is the variance of the shape samples o^s for that voxel. Unfortunately, there often exist several voxels with the same variance and choosing one to explore is non-trivial. However, we found that most uncertain voxels form small clusters. We chose to explore the cluster with the most flat surface, which is advantageous for making a robust contact with the object.

B. Sampling of Shapes

One of the crucial parts in Section III-A is the sampling of shapes from the probability distribution $P(O_t|v, h_{1:t-1})$. Previous work on probabilistic shape completion [12] achieved this by training a 3D CNN to reconstruct voxelized objects and using the variational inference technique Monte Carlo dropout [25] for sampling. However, for this process to work on visuo-haptic data, it requires training the CNN on both haptic and visual data, with haptic data collected from random positions all around the object. Unfortunately, no such dataset exists, and curating one is expensive because of the haptic data collection process [26].

Instead, we propose to train a reconstruction network that can accurately reconstruct shapes based on visuo-haptic data without explicitly training on such data. For this, we chose to use the IGR architecture [4] that learns the SDF of the underlying surface. IGR is a Multi-Layer Perceptron (MLP) $f(\mathbf{x}; \boldsymbol{\theta}, \cdot): \mathbb{R}^3 \rightarrow \mathbb{R}$, where \mathbf{x} is a 3D point and the parameters $\boldsymbol{\theta}$ are trained such that f is approximately the SDF to a plausible surface \mathcal{M} defined by the point cloud $\mathcal{X} = \{\mathbf{x}_c\}_{c \in C}$ and optionally the point normals $\mathcal{N} = \{\mathbf{n}_c\}_{c \in C}$, where C is the set of points in the point cloud. The \cdot is an additional parameter which is introduced below. The loss function to train IGR is

$$\ell_{rec}(\boldsymbol{\theta}, \cdot) = \ell_{\mathcal{X}}(\boldsymbol{\theta}, \cdot) + \lambda \mathbb{E}_{\mathbf{x}} [\|\nabla_{\mathbf{x}} f(\mathbf{x}; \boldsymbol{\theta}, \cdot)\| - 1]^2, \quad (4)$$

where $\lambda > 0$, and

$$\ell_{\mathcal{X}}(\boldsymbol{\theta}, \cdot) = \frac{1}{|C|} \sum_{c \in C} (|f(\mathbf{x}_c; \boldsymbol{\theta}, \cdot)| + \tau \|\nabla_{\mathbf{x}} f(\mathbf{x}_c; \boldsymbol{\theta}, \cdot) - \mathbf{n}_c\|). \quad (5)$$

The first term in Eq. 4, which is detailed in Eq. 5, pushes f to vanish on \mathcal{X} . If normal data exist (our case), then $\tau := 1$ and $\nabla_{\mathbf{x}} f$ is pushed to the supplied normals \mathcal{N} . The second term in Eq. 4, called the Eikonal term, regularizes the network to produce smooth reconstructions by forcing the gradients of $\nabla_{\mathbf{x}} f$ to be of unit 2-norm.

By default, a separate IGR is trained for every single shape. However, this prohibits sampling from the underlying shape distribution $P(O_t|v, h_{1:t-1})$. Therefore, we chose to train a multi-shape IGR, which is realized by first selecting a separate latent vector \mathbf{z}_j for each training example $j \in J$ and then train the network to approximate the SDF associated with each \mathbf{z}_j . The multi-shape IGR takes the following form $f(\mathbf{x}; \boldsymbol{\theta}, \mathbf{z}_j)$.

To complete the shape of an object from a partial point cloud with the multi-shape IGR comes down to finding a latent code $\hat{\mathbf{z}}_i$, where $i \in I$ are test samples, that best reconstructs the SDF of the point cloud. To find such a latent code, we treat the observed point cloud as the ground truth and use

gradient optimization to fine-tune an initially random code $\hat{z}_{i,0}$. Formally, this fine-tuning is expressed as

$$\hat{z}_{i,t} = \hat{z}_{i,t-1} - \alpha \nabla_{\hat{z}_{i,t-1}} \ell(\theta, \hat{z}_{i,t-1}), \quad (6)$$

where α is the step-size and $\nabla_{\hat{z}_{i,t-1}}$ is the gradient of the following loss function:

$$\ell(\theta, \hat{z}_{i,t-1}) = \ell_{rec}(\theta, \hat{z}_{i,t-1}) + \gamma \|\mathbf{z}_{i,t-1}\|. \quad (7)$$

We chose $\gamma = 0.01$, and $\ell_{rec}(\theta, \hat{z}_{i,t-1})$ is the loss in Eq. 4.

For drawing shape samples from a multi-shape IGR, two alternatives exist. The first, and most obvious, is to sample multiple latent codes $\hat{\mathbf{z}}_{1:S}$, where S is the number of latent codes sampled, and optimize each of them individually for a fixed number of gradient descent steps. The second option, which we chose to use in our experiments, is to sample and optimize only one latent code \hat{z}_1 and select S intermediate optimized codes as the shape samples. For instance, if we optimized the latent code for 800 steps, we could select the latent code after 650, 700, 750 and 800 steps as our samples. This resembles Metropolis sampling in that samples are generated from a supposedly converged Markov chain, however, in our case we do not use the Metropolis rejection rule in the optimization process for simplicity but the stochasticity is introduced through sampling mini-batches in the optimization. This option is significantly faster than the first one and was empirically found to provide similar results.

C. Active Visuo-Haptic Object Shape Completion

We propose Algorithm 1 for active shape completion by combining the probabilistic shape completion from Section III-A with the sampling of shapes in Section III-B (see also Fig. 1). The algorithm starts by generating a random latent code \hat{z}_0 (line 5). Then that latent code is optimized with gradient descent over the current point cloud and intermediate codes \hat{z}_g are saved (lines 8–14). From the intermediate latent codes, meshes are reconstructed and transformed to voxel grids for the variance computation (line 15). Next, the voxel to touch (as described in Section III-A) is computed, explored, and the information is added to the point cloud (lines 16–18). Note that if no collision is detected at the target location (the robot is commanded to move on a straight line towards the target and 10 cm beyond), no point cloud is saved and the robot returns to the start position and selects a new position for exploration. After all M haptic explorations are done, the last latent code is optimized and the final shape is reconstructed (lines 19–23). Some steps are illustrated in the accompanying video at <https://youtu.be/iZF4ph4zMEA>.

D. Implementation details

The IGR network was implemented in PyTorch 1.0.0. The network structure was the same as in [4] and consisted of 8 fully connected layers with 512 neurons and a skip connection in the 4th layer. The training was carried out on NVIDIA GeForce GTX 1080 Ti for 3500 iterations, with a batch size of 8 and a latent vector size of 256.

To train the network, we curated our own dataset of 87 unique meshes from the YCB [27] and Grasp Database [28]

Algorithm 1 Active Visuo-Haptic Shape Completion

- 1: **Inputs:** point cloud \mathbf{P} , number of haptic explorations M , number of gradient-descent steps G , steps before storing latent shape L
- 2: **Output:** Final shape completion O
- 3: $\mathbf{H} \leftarrow \emptyset$ ▷ Empty set of haptic data
- 4: $\mathbf{P}_0 \leftarrow \mathbf{P}$
- 5: $\hat{z}_0 \leftarrow$ Sample initial latent code
- 6: **for** $m \leftarrow 1, \dots, M$ **do**
- 7: $\mathbf{Z} \leftarrow \emptyset$ ▷ Empty set of latent codes
- 8: **for** $g \leftarrow 1, \dots, G$ **do**
- 9: $\hat{z}_g \leftarrow$ Optimize \hat{z}_{g-1} over \mathbf{P}_{m-1} using Eq. 6
- 10: **if** $g \bmod L == 0$ **then**
- 11: $\mathbf{Z} \leftarrow \mathbf{Z} + \hat{z}_g$
- 12: **end if**
- 13: **end for**
- 14: $\hat{z}_0 \leftarrow \hat{z}_g$
- 15: $\mathbf{V} \leftarrow$ Reconstruct shapes from \mathbf{Z} and calculate their variance
- 16: $h_m \leftarrow$ Calculate next touch using Eq. 3 and \mathbf{V}
- 17: $\mathbf{H} \leftarrow$ Execute h_m and append the touch point
- 18: $\mathbf{P}_m \leftarrow \mathbf{P} + \mathbf{H}$
- 19: **end for**
- 20: **for** $g \leftarrow 1, \dots, G$ **do**
- 21: $\hat{z}_g \leftarrow$ Optimize \hat{z}_{g-1} over \mathbf{P}_I using Eq. 6
- 22: **end for**
- 23: $O \leftarrow$ Reconstruct shape using \hat{z}_g

datasets. Each mesh was centered at the origin and scaled such that the longest dimension was between -1 and 1. To generate the ground truth point cloud of a mesh, we sampled 100000 points evenly over the complete object and, for each point, also estimated its normal. We rotated each mesh into 16 different views, resulting in 1392 training samples in total. Using the same procedure, we also generated a test set of 35 completely novel objects from both datasets.

IV. EXPERIMENTS

The experiments address the following two questions:

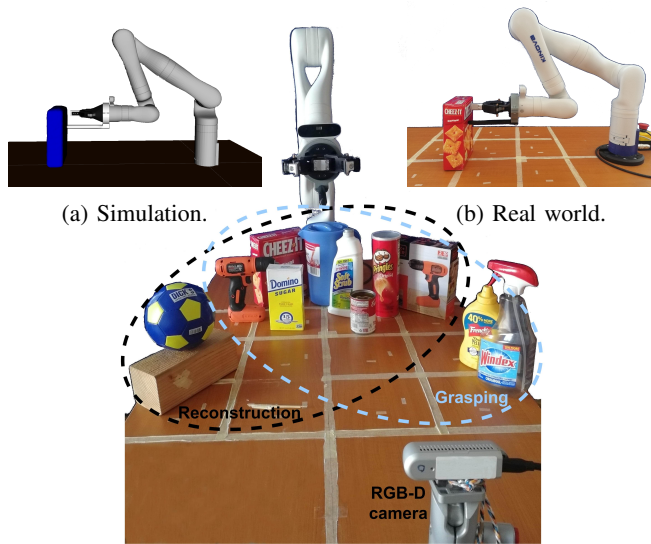
- 1) What is the shape reconstruction accuracy of Act-VH?
- 2) What is the impact of Act-VH on grasp success rate?

To reliably answer these questions, we conducted two experiments. The first experiment (Section IV-B) evaluated shape reconstruction in simulation and in the real world, while the second experiment (Section IV-C) evaluated grasp success rates on real hardware.

A. Experimental setup

In both the simulation and real-world experiments, we used a Kinova Gen3 robot equipped with a custom made finger to perform haptic exploration and a Robotiq 2F-85 gripper for grasping. In the real-world experiments, we used an Intel RealSense D430 depth camera to capture the point cloud – see Fig. 2.

For haptic exploration, we moved the robot along a pre-defined approach direction until the torques of the robot joints



(c) The real-world setup with the Kinova Gen3 robotic arm, the RGB-D camera (always in front of the robot), and the objects used for reconstruction and for grasping.

Fig. 2: Simulated (a) and real (b,c) environment.

crossed a pre-defined threshold. The global position of the finger was then transformed into the same reference system as the visual point cloud. The objects were attached to the table using a double-sided tape. If nothing specific is noted, we maximally executed five touches in all experiments to keep the total execution time low. The time taken to run the pipeline with five touches is about 5 minutes on average.

In the reconstruction experiments, we benchmarked Act-VH against five other baseline methods: Ball Pivoting Algorithm (BPA) [29], Poisson reconstruction (Poisson) [30], Convex Hull reconstruction (Hull), Alpha shapes (Alpha) [31], and GPIS [32]. BPA [29] reconstructs a shape by rolling a sphere with a pre-defined radius over all points, and if three points are inside the sphere, they are connected with a triangle. Poisson reconstruction [30] solves an optimization problem that creates a smooth surface over the points but can only reconstruct the visible part of the surface. The Hull method reconstructs the input point cloud with a convex hull over the points. Alpha is a generalized Hull method that smooths the object surface and can also remove volume from the inside of concave objects. GPIS [32] trains a GP to reconstruct the implicit surface of the object. For GPIS, we used similar hyper-parameters as reported in [22].

To benchmark other methods, we still used Act-VH to calculate the reconstruction uncertainty and where to touch the object but used a baseline method instead of IGR for the final reconstruction. In the reconstruction experiments, we evaluated all methods using the Chamfer distance and Jaccard similarity, while in the grasping experiment, we used the grasp success rate. For calculating the Jaccard similarity, each mesh was voxelized into a voxel grid of size 40^3 .

B. Object Reconstruction

In this experiment, we evaluated the reconstruction accuracy in simulation and in the real world. To evaluate in simulation,

we developed our own visuo-haptic robotic simulation environment in the MuJoCo physics simulator [33]. The environment, shown in Fig. 2a, consist of a robot, an object mesh, and a virtual camera for capturing the point cloud of the object. In both simulation and the real world, the robot planned and moved to the location we wanted to haptically explore and stopped once contact was detected as shown in Fig. 2b.

In simulation, we evaluated the reconstruction accuracy on 35 test objects, while in the real world, we used the 10 objects shown in Fig. 2c that were selected because they differed in size and shape. Each reconstruction was repeated three times for each object and method combination, resulting in 105 unique reconstructions per method in simulation and 30 in the real world. In the simulation experiment, we further compared Act-VH to: (i) a random policy that touched the first reachable voxel from the set of uniformly sampled voxels on the surface of the reconstruction, and (ii) a GPIS-driven policy, where the voxel with the largest standard deviation was selected [17].

Note that GPIS approximates the surface covered by the input points but does not assume a closed volume. Therefore, we needed to select a reasonable first touch point heuristically.

Fig. 3 shows the reconstruction results separately for simulation (Fig. 3a) and real world (Fig. 3b). Overall, the reconstruction accuracy for all methods improved with the number of touches, meaning that Jaccard similarity increased and Chamfer distance decreased. Furthermore, both simulation and real world results follow the same trend indicating the robustness of our method.

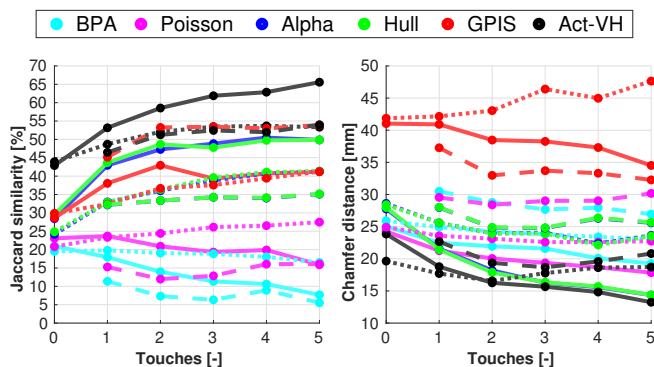
Based on the results in Fig. 3, we can clearly see that Act-VH outperforms all other baselines across the board. For example, Fig. 3a shows that Act-VH outperforms the random one already after one touch, and after five touches Act-VH reaches around 10-20% higher Jaccard similarity than the random one and over 5 mm lower Chamfer distance.

The results in Fig. 3a also show that Act-VH exploration outperforms GPIS exploration. The reason GPIS performs poorly is because it requires touches to be evenly distributed around the whole object.¹ Other reconstruction works presented similar results [17], [22], where in the haptic-only work [17], more than 100 touches were required to reconstruct the front side only, while in the visuo-haptic work [22], more than 20 touches were collected and the results were still poor.

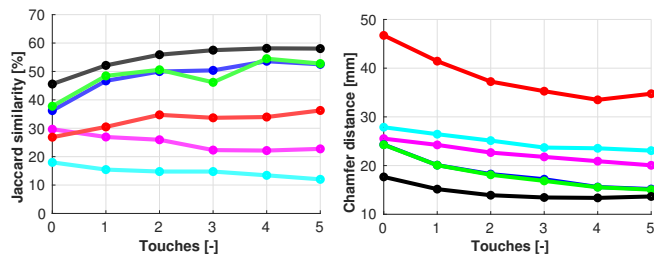
When comparing reconstruction methods, the second best method was Hull, with a Chamfer distance around 15 mm. One reason Hull achieved such a low Chamfer distance is that it creates more sharp reconstructions, which strongly influences the Chamfer distance. However, for Hull to produce good reconstructions the point cloud must contain points covering the whole object, which is in practice best achieved with Act-VH exploration.

Interestingly, with Act-VH exploration, the Jaccard similarity of BPA and Poisson decreased with more touches, while with a random policy, it stayed flat or increased slightly. These results point to the fact that BPA and Poisson can only reconstruct the visible part of the surface, *i.e.*, parts

¹Note that we used the exponential kernel. Results with other kernels (*e.g.*, thin plate kernel as in [3]) may be different.



(a) Simulation results. Solid lines represent the accuracy obtained when executing Act-VH touches, dotted lines random touches, and dashed lines GPIS-uncertainty-driven touches.



(b) Real world results.

Fig. 3: The average reconstruction accuracy from simulation (a) and the real world (b). For Jaccard similarity, larger values are better; for Chamfer distance, smaller values are better.

that are covered by the point cloud, and random exploration has a high chance of being close to those points resulting in more useful information. In contrast, Act-VH exploration most likely returns a data point far from the visual point cloud, leading BPA and Poisson to create strange artifacts that, once voxelized, result in lower Jaccard similarity. Although the Act-VH policy resulted in lower Jaccard similarity for BPA and Poisson, it still outperformed the random policy in terms of Chamfer distance.

Fig. 4 shows example reconstructions after 5 touches. These examples show that: (i) BPA and Poisson are unable to complete the whole object accurately, (ii) Alpha and Hull reconstruct very sharp and unrealistic objects, and (iii) GPIS is poor at reconstructing the object where no points are available. In contrast, Act-VH can capture both global and local features, resulting in smooth and faithful reconstructions. A challenging object to shape complete was the partly transparent spray bottle in the bottom row of Fig. 4. Nevertheless, Act-VH still reconstructed it quite well compared to the ground truth and other methods. The results of an incremental Act-VH reconstruction with five touches is visualized in Fig. 5, highlighting that if the first reconstruction is good, which happened to be the case in simulation, additional touches only locally refine the objects. However, if the initial reconstruction is poor, which was the case in the real world experiment, additional touches lead to more global refinements. Note that the initial estimation could be improved by replacing random sampling for mini-batches with Farthest Point Sampling (FPS) (as in [34]) which better preserves the global information about the object. However, our experiments showed that after haptic

exploration, FPS is not leveraging this information well and is outperformed by random selection.

Finally, we investigated if the reconstruction accuracies in Fig. 3 approach some steady-state value with more touches. We let Act-VH explore three objects three times in simulation with 50 touches. The results are presented in Fig. 6. Based on these results, it seems that Act-VH does approach a steady-state Chamfer distance after about 20 touches, albeit some fluctuations are still present. The same conclusion cannot be made for the Jaccard similarity, which actually gets worse after about 25 touches. The primary reason the Jaccard similarity starts to decrease was due to errors in the exact location of contact which originated from imprecise joint torque collision detection. Although Act-VH can cope with some errors, ultimately, after enough touches, the performance starts to decrease.

C. Robotic Grasping

The final experiment evaluates the impact of active visuo-haptic object shape completion on grasp success rate. We used the same overall setup as in the real-world reconstruction experiment, but changed two of the objects. The ten objects we used are shown in Fig. 2c. All of these objects, except the yellow mustard bottle, were completely new. We decided to benchmark Act-VH against the Hull method because it reached the second-best reconstruction accuracy on simulated and real objects.

For planning grasps on the reconstructed objects, we used the simulated annealing planner in GraspIt! that ran for 75000 steps. Out of the planned grasps, the first physically reachable grasp with the highest ϵ -quality metric was executed on the robot. To study the effect of the number of touches on grasp success rate, we planned and executed a grasp after zero, three, and five touches. We repeated the reconstruction and grasping procedure three times for each combination of objects, method, and the number of touches, resulting in 180 grasps in total. The robot performed a grasp by first picking the object, then moving 10 cm upwards, and finally rotating the last joint $\pm 90^\circ$. The grasp was considered successful if the robot did not drop the object during this movement; otherwise, it was unsuccessful.

Fig. 7a shows the average grasp success rates over varying number of touches. We can clearly see that Act-VH is superior to Hull. For instance, after five touches, Act-VH achieved an 80% average grasp success rate while Hull only achieved 46.7%. As expected, the grasp success rate with Act-VH improves with the number of touches, from 38% to 80%. In comparison, the success rate for Hull was unchanged between 3 and 5 touches, indicating that the additional haptic data did not improve the convex hull reconstruction for grasp planning. This fact is highlighted in Fig. 4, where the Hull reconstruction

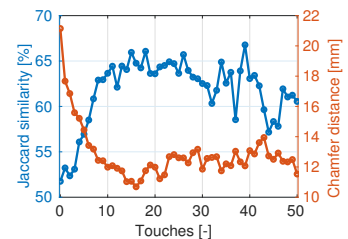


Fig. 6: Long exploration – 50 touches in simulation. Average Jaccard similarity and Chamfer distance for three objects over three repetitions.

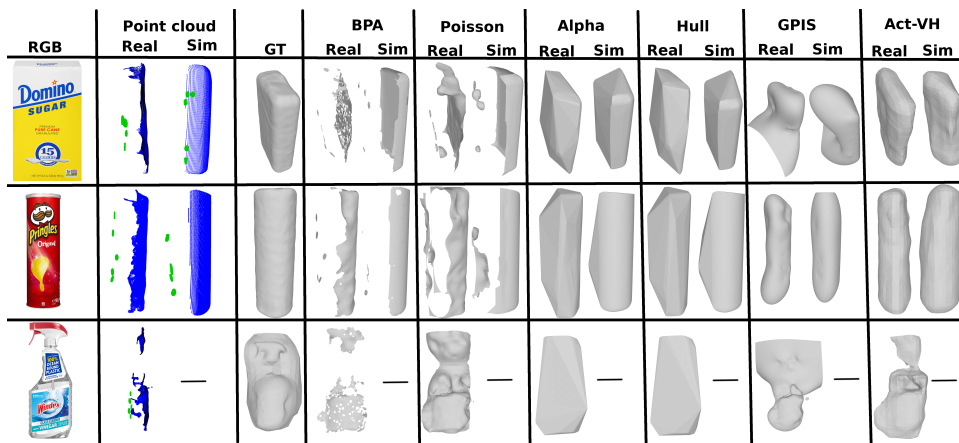


Fig. 4: Reconstruction examples in simulation (Sim) and real with all methods on three objects: A basic rectangular object (upper row), the object for which Act-VH achieved the worst grasp success rate (middle row), and an adversarial object (bottom row). Touches in the point clouds are highlighted with green color. There is no reconstruction in simulation for the adversarial object because there was no appropriate ground truth mesh.

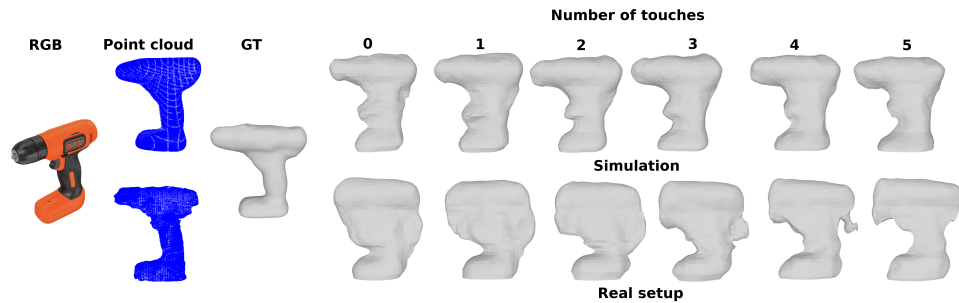


Fig. 5: An example reconstruction of an object after each touch with Act-VH in simulation (upper row) and the real world (bottom row).

creates “slopes” from one set of points to another which the grasp planner seemed to favor but resulted in unsuccessful grasps. The same reconstruction artifacts were not visible for Act-VH. Despite these differences, haptic exploration still increased the grasp success rate by more than 100% between zero and five touches irrespective of the methods, highlighting the benefit of better reconstructions.

Finally, Fig. 7b shows the average grasp success rates of Act-VH and Hull after five touches on the ten objects individually. The results indicate that Act-VH performs better than or on par with Hull on all objects. Hull has particular problems with larger objects, such as objects 1, 2, 4, 5, and 10, which most probably stem from the “slope” artifacts mentioned earlier. Act-VH, on the other hand, only performs poorly on object 7, where the two failed grasps were side-grasps for which the object slipped out of the gripper because of its short fingers. One possible reason why more side grasps were produced was because the reconstructed object was much thinner than in the real world (shown in the center row of Fig. 4).

V. CONCLUSIONS AND FUTURE WORK

We presented Act-VH, an active visuo-haptic shape completion method. The challenge in visuo-haptic shape completion is to decide the most informative touch location. To this end, Act-VH uses a probabilistic shape completion network to assess where the reconstruction is most uncertain. This location is then used for haptic exploration. The reconstruction accuracy

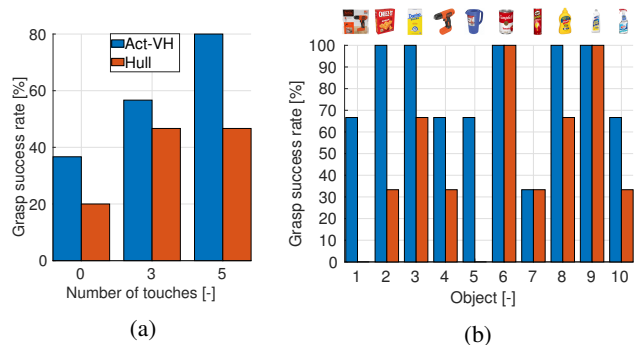


Fig. 7: Average grasp success rate of IGR (blue) and Hull (red).

of Act-VH compared to five baseline methods shows, both in simulation and real world, that Act-VH produces the best reconstructions and that its reconstruction accuracy increases most with the number of haptic explorations. Furthermore, active visuo-haptic shape completion was also beneficial for robotic grasping where Act-VH reached significantly higher grasp success rates than the Hull method.

To assess the uncertainty of current shape reconstruction, we sampled latent codes during shape optimization and used the variance of the reconstructed shapes from this phase as a measure of uncertainty. Note that this method may not reflect the true uncertainty about the object shape given the available information. An alternative method—sampling and

then optimizing multiple latent codes independently—did not yield better results and was computationally more expensive. However, this remains an empirical result and better theoretical grounding would be required.

Although Act-VH achieved promising reconstruction accuracy and grasp success rates, there is still room for improvements. One improvement is to modify the loss function of IGR to also incorporate data points that we know are not on the surface. For instance, haptic exploration does not only indicate where the surface exists, but also where it does not exist. Another improvement is to also model the haptic location as uncertain. This would allow to select touch locations that are most robust to shape and robot uncertainties. Practically, performance would increase with a more sensitive contact detection method using e.g. a F/T sensor in the robot wrist or tactile sensors at the fingertip.

In summary, the work presented here shows that we can achieve accurate shape reconstructions with active visuo-haptic shape completion. This, in turn, enables the use of visuo-haptic exploration in perceptually uncertain environments such as cluttered scenes where objects occlude each other.

REFERENCES

- [1] J. Varley, C. DeChant, A. Richardson, J. Ruales, and P. Allen, "Shape completion enabled robotic grasping," in *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, Sept. 2017, pp. 2442–2447.
- [2] S. Wang, J. Wu, X. Sun, W. Yuan, W. T. Freeman, J. B. Tenenbaum, and E. H. Adelson, "3D Shape Perception from Monocular Vision, Touch, and Shape Priors," in *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, Oct. 2018, pp. 1606–1613.
- [3] M. Björkman, Y. Bekiroglu, V. Högman, and D. Kragic, "Enhancing visual perception of shape through tactile glances," in *2013 IEEE/RSJ International Conference on Intelligent Robots and Systems*, Nov. 2013, pp. 3180–3186.
- [4] A. Gropp, L. Yariv, N. Haim, M. Atzmon, and Y. Lipman, "Implicit Geometric Regularization for Learning Shapes," in *International Conference on Machine Learning*. PMLR, Nov. 2020, pp. 3789–3799.
- [5] R. Bajcsy, Y. Aloimonos, and J. K. Tsotsos, "Revisiting active perception," *Autonomous Robots*, vol. 42, no. 2, pp. 177–196, 2018.
- [6] J. Bohg, M. Johnson-Roberson, B. León, J. Felip, X. Gratal, N. Bergström, D. Kragic, and A. Morales, "Mind the gap-robotic grasping under incomplete observation," in *2011 IEEE International Conference on Robotics and Automation*. IEEE, 2011, pp. 686–693.
- [7] R. Schnabel, P. Degener, and R. Klein, "Completion and reconstruction with primitive shapes," in *Computer Graphics Forum*, vol. 28, no. 2. Wiley Online Library, 2009, pp. 503–512.
- [8] M. Pauly, N. J. Mitra, J. Giesen, M. H. Gross, and L. J. Guibas, "Example-based 3d scan completion," in *Symposium on Geometry Processing*, no. CONF, 2005, pp. 23–32.
- [9] M. Li, K. Hang, D. Kragic, and A. Billard, "Dexterous grasping under shape uncertainty," *Robotics and Autonomous Systems*, vol. 75, pp. 352–364, 2016.
- [10] A. Dai, C. R. Qi, and M. NieBner, "Shape Completion Using 3D-Encoder-Predictor CNNs and Shape Synthesis," in *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. Honolulu, HI: IEEE, July 2017, pp. 6545–6554.
- [11] X. Han, Z. Li, H. Huang, E. Kalogerakis, and Y. Yu, "High-Resolution Shape Completion Using Deep Neural Networks for Global Structure and Local Geometry Inference," in *2017 IEEE International Conference on Computer Vision (ICCV)*. Venice: IEEE, Oct. 2017, pp. 85–93.
- [12] J. Lundell, F. Verdoja, and V. Kyrki, "Robust Grasp Planning Over Uncertain Shape Completions," in *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, Nov. 2019, pp. 1526–1532.
- [13] J. J. Park, P. Florence, J. Straub, R. Newcombe, and S. Lovegrove, "DeepSDF: Learning Continuous Signed Distance Functions for Shape Representation," in *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. Long Beach, CA, USA: IEEE, June 2019, pp. 165–174.
- [14] M. Atzmon and Y. Lipman, "SAL: Sign Agnostic Learning of Shapes From Raw Data," in *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. Seattle, WA, USA: IEEE, June 2020, pp. 2562–2571.
- [15] Z. Wu, S. Song, A. Khosla, F. Yu, L. Zhang, X. Tang, and J. Xiao, "3D ShapeNets: A deep representation for volumetric shapes," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 1912–1920.
- [16] S. Ottenhaus, M. Miller, D. Schiebener, N. Vahrenkamp, and T. Asfour, "Local implicit surface estimation for haptic exploration," in *2016 IEEE-RAS 16th International Conference on Humanoid Robots (Humanoids)*, Nov. 2016, pp. 850–856.
- [17] Z. Yi, R. Calandra, F. Veiga, H. van Hoof, T. Hermans, Y. Zhang, and J. Peters, "Active tactile object exploration with Gaussian processes," in *2016 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, Oct. 2016, pp. 4925–4930.
- [18] D. Driess, P. Englert, and M. Toussaint, "Active learning with query paths for tactile object shape exploration," in *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, Sept. 2017, pp. 65–72.
- [19] S. Dragiev, M. Toussaint, and M. Gienger, "Uncertainty aware grasping and tactile exploration," in *2013 IEEE International conference on robotics and automation*. IEEE, 2013, pp. 113–119.
- [20] G. Z. Gandler, C. H. Ek, M. Björkman, R. Stolkin, and Y. Bekiroglu, "Object shape estimation and modeling, based on sparse Gaussian process implicit surfaces, combining visual data and tactile exploration," *Robotics and Autonomous Systems*, vol. 126, p. 103433, Apr. 2020.
- [21] S. Ottenhaus, D. Renninghoff, R. Grimm, F. Ferreira, and T. Asfour, "Visuo-Haptic Grasping of Unknown Objects based on Gaussian Process Implicit Surfaces and Deep Learning," in *2019 IEEE-RAS 19th International Conference on Humanoid Robots (Humanoids)*, Oct. 2019, pp. 402–409.
- [22] D. Watkins-Valls, J. Varley, and P. Allen, "Multi-Modal Geometric Learning for Grasping and Manipulation," in *2019 International Conference on Robotics and Automation (ICRA)*, May 2019, pp. 7339–7345.
- [23] E. Smith, R. Calandra, A. Romero, G. Gkioxari, D. Meger, J. Malik, and M. Drozdal, "3D Shape Reconstruction from Vision and Touch," in *Advances in Neural Information Processing Systems*, vol. 33. Curran Associates, Inc., 2020, pp. 14 193–14 206.
- [24] E. Smith, D. Meger, L. Pineda, R. Calandra, J. Malik, A. Romero Soriano, and M. Drozdal, "Active 3D Shape Reconstruction from Vision and Touch," *Advances in Neural Information Processing Systems*, vol. 34, 2021.
- [25] Y. Gal and Z. Ghahramani, "Dropout as a bayesian approximation: Representing model uncertainty in deep learning," in *international conference on machine learning*. PMLR, 2016, pp. 1050–1059.
- [26] S. Luo, J. Bimbo, R. Dahiya, and H. Liu, "Robotic tactile perception of object properties: A review," *Mechatronics*, vol. 48, pp. 54–67, 2017.
- [27] B. Calli, A. Singh, A. Walsman, S. Srinivasa, P. Abbeel, and A. M. Dollar, "The ycb object and model set: Towards common benchmarks for manipulation research," in *2015 International Conference on Advanced Robotics (ICAR)*, 2015, pp. 510–517.
- [28] D. Kappler, J. Bohg, and S. Schaal, "Leveraging big data for grasp planning," in *2015 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2015, pp. 4304–4311.
- [29] F. Bernardini, J. Mittleman, H. Rushmeier, C. Silva, and G. Taubin, "The ball-pivoting algorithm for surface reconstruction," *IEEE Transactions on Visualization and Computer Graphics*, vol. 5, no. 4, pp. 349–359, Oct. 1999.
- [30] M. Kazhdan, M. Bolitho, and H. Hoppe, "Poisson surface reconstruction," in *Proceedings of the fourth Eurographics symposium on Geometry processing*, vol. 7, 2006.
- [31] B. Guo, J. Menon, and B. Willette, "Surface Reconstruction Using Alpha Shapes," *Computer Graphics Forum*, vol. 16, no. 4, pp. 177–190, 1997.
- [32] O. Williams and A. Fitzgibbon, "Gaussian process implicit surfaces," in *Gaussian Processes in Practice*, 2006.
- [33] E. Todorov, T. Erez, and Y. Tassa, "Mujoco: A physics engine for model-based control," in *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2012, pp. 5026–5033.
- [34] C. R. Qi, L. Yi, H. Su, and L. J. Guibas, "Pointnet++: Deep hierarchical feature learning on point sets in a metric space," in *Proceedings of the 31st International Conference on Neural Information Processing Systems*, ser. NIPS'17. Red Hook, NY, USA: Curran Associates Inc., 2017, p. 5105–5114.