

CLASP: Constrained Latent Shape Projection for Refining Object Shape from Robot Contact

Anonymous Author(s)

Affiliation

Address

email

1 **Abstract:** Robots need both visual and contact sensing to effectively estimate
2 the state of their environment. Camera RGBD data provides rich information of
3 the objects surrounding the robot, and shape priors can help correct noise and
4 fill in gaps and occluded regions. However, when the robot senses unexpected
5 contact, the estimate should be updated to explain the contact. To address this
6 need, we propose CLASP: Constrained Latent Shape Projection. This approach
7 consists of a shape completion network that generates a prior from RGBD data
8 and a procedure to generate shapes consistent with both the network prior and
9 robot contact observations. We find CLASP consistently decreases the Chamfer
10 Distance between the predicted and ground truth scenes, while other approaches
11 do not benefit from contact information.

12 1 Introduction

13 You look into a cabinet and see a box of crackers. You reach in and attempt to grab the box from
14 the side, but your fingers hit something. Perhaps this box is larger than you thought? Your mental
15 model of the box updates, you try a wider grasp, and you successfully retrieve your snack. Robots are
16 currently not so adept. While they can estimate the pose of known shapes [1] or estimate parameters
17 of objects [2], they cannot yet fuse this visual and contact information to draw from the wide range
18 of shape priors in the world. A robot could try to learn its next action directly from vision and force
19 feedback instead [3], but this approach lacks the logic to generalize to scenarios not seen in training.

20 We propose a method that allows robots to mimic the process of updating object shape from contact
21 information. A shape completion neural network first generates beliefs over possible object shapes
22 based on visual RGBD data. The belief updates the object shape to be consistent with contact
23 information gathered by a robot moving in the scene. We make the realistic assumption that the
24 RGBD camera perceiving the scene suffers from sensor noise and occlusion. We assume the robot
25 can sense *if* it collides with an object, but not *where* the contact was made (i.e., no sensorized skin).
26 Many of the “cobot” platforms available today utilize this contact model to detect collision and stop
27 before harming a person.

28 Formally, this type of contact creates a contact manifold, a thin space of shapes with a boundary
29 bordering the robot. Past work has projected shapes onto the contact manifold in object pose space
30 [4] and robot configuration space [1], but both require known shape geometry. Our objective is to
31 update the unknown shape geometry to satisfy contact constraints. Returning to the cracker box
32 example, the robot will be unsure if the contact occurred at the top finger, the bottom finger, or
33 perhaps the back of the hand or the elbow. Filling in all possible contact points would lead to absurd
34 scenes with robot shells protruding from the cracker box. However, ignoring this contact information
35 leaves the robot with the original belief of the thinner cracker box and no explanation of why the
36 attempted grasp failed. Our shape completion network generates a prior in latent shape space which
37 can be decoded into shapes in workspace; however, shapes generated directly from this latent prior
38 are unlikely to satisfy contact constraints.

39 Our *key insight* is that latent samples from our neural network can be projected onto the contact
40 manifold in the latent shape space using iterative gradient descent, creating shapes both likely under

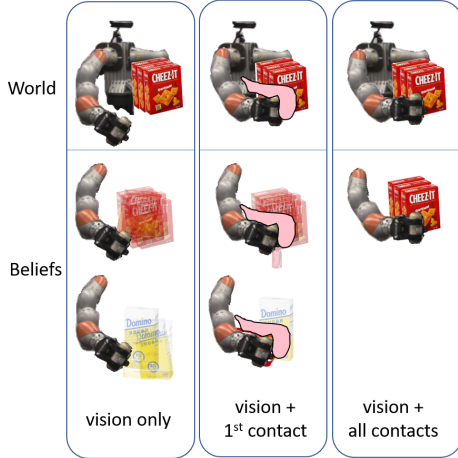


Figure 1: A visual RGBD view of objects leave ambiguity final shape due to sensor noise and occlusion, which we store as a set of sampled scenes in a particle filter. Contact information (pink) reduces ambiguity, and using CLASP the particles converge to the true shape.

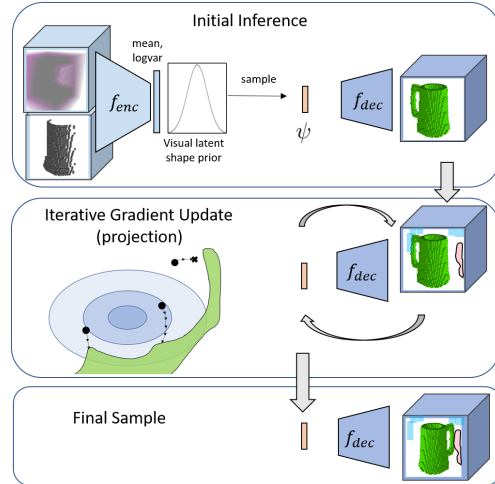


Figure 2: CLASP Architecture

Top: Shapes sampled from RGBD using PSSNet [5].
Middle Right: A robot motion detects free space (light blue) and a collision set (pink).
Middle Left: Latent samples are projected to satisfy contacts (green). Ovals depict the latent prior.
Bottom: Final samples satisfy the contact constraints.

41 the visual prior and consistent with the contact information. We further expect these projected shapes
 42 to be closer to ground truth than direct samples not considering contact.

43 We accomplish this with our proposed Constrained LATent Shape Projection (CLASP), which stores
 44 a belief over shapes in a particle filter. Each particle represents a collection of latent object shapes
 45 which can be decoded into a scene. Every new robot measurement of contact and freespace triggers
 46 an update on all particles. During each particle update, gradient steps are taken to increase the
 47 occupancy likelihood of the most likely point(s) explaining the contact(s), decrease the occupancy
 48 likelihood of the free points, and increase the latent likelihood under the shape prior.

49 We test this method both in simulation and on a live robot by constructing scenes of objects on a
 50 tabletop, generating robot motions that generate freespace and contact observations (Fig. 1), updat-
 51 ing the belief using CLASP as well as baselines, and comparing the set of sampled shapes to the
 52 ground truth scene. We find CLASP outperforms both ablations of CLASP, as well as the approaches
 53 of Rejection Sampling and directly updating the input to the shape completion network. We also
 54 find that CLASP produces more accurate scenes than a VAE.GAN shape completion network.

55 2 Related Work

56 **Shape Completion from Vision:** The goal of Shape Completion is to predict a full shape from a
 57 single partial input. Recently, neural networks have become a popular method of shape completion.
 58 A common network architecture learns an encoder to a feature space followed by a decoder to the
 59 shape output [6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17]. Scene Completion networks are trained
 60 on larger spatial volumes of occupied points and use similar architectures with adaptations to join
 61 information from multiple scales [18]. While most methods predict the single best estimated shape,
 62 we build off work that uses a variational autoencoder architecture to produce plausible and diverse
 63 shapes [5]. We draw from the vast work on shape and scene completion and contribute a method
 64 that improves scene estimates using contact information from a robot.

65 **Shape Completion from Touch:** In different works, “touch” can refer to a single known contact
 66 point, a contact configuration, force-torque measurements, or a rich tactile sensor. Work using the
 67 definition of contact point or contact configuration typically uses touch to reduce the version space
 68 of shape possibilities [19], but such approaches cannot tractably capture the diversity of all shapes.

69 Alternatively, some situations model known shapes with unknown poses [2]. Both classical Iterative
 70 Closest Point [20, 21] and neural networks [22, 23, 24] have been used to predict valid poses. To
 71 generate samples consistent with contact information, the Implicit Manifold Particle Filter projects
 72 sampled poses onto the contact manifold using an iterative approach [4]. Analogously, we project
 73 sampled particles onto the contact manifold in the latent shape space of our neural network.

74 Touch can also refer to the rich tactile sensors such as GelSight [25] or soft-bubble grippers [26],
 75 with input more analogous to images. Tactile patches can be directly mapped to visual features [27].
 76 Alternatively, neural networks have used these sensors for material classification [28] and grasped
 77 pose estimation [29]. We do not assume our contact sensing has such rich information.

78 **Combining Vision and Touch:** A neural network can combine vision and touch (force + torque
 79 [3], or GelSight [30]) using separate encoders to a latent space for each sensing modality alongside a
 80 decoder to a variety of spaces. We considered a similar encoder structure with a decoder to produce
 81 completed shapes, but this would require a large dataset of (Shapes \times Contact + Freespace) mea-
 82 surements and the resulting network would be only applicable to the robot used for training. The loss
 83 of a neural network can be tweaked during training to bias towards priors similar to contact, such as
 84 connectivity and stability [31]. A neural network can output a downstream objective, such as grasp
 85 success probability, instead of shape reconstruction, and thus may be successful for a diverse array
 86 of objects where accurate reconstruction is not available [32]. Instead of a neural network, surface
 87 reconstruction has been done using a Gaussian Process (GP) prior fit to tactile measurements [33, 34].

88 Our method is most similar to the work of Wang et al. [35], which uses gradient descent on the
 89 latent space of a shape completion network to enforce touch constraints. Where that work uses a
 90 high-resolution GelSight tactile sensor to refine shape details previously reconstructed from vision,
 91 our work focuses on ambiguous shapes (e.g. a box with unknown depth, or novel shapes not in the
 92 training data) and the lower information measurement of contact detection. We accomplish much
 93 larger shape updates by using a diverse set of predictions and a novel projection loss function.

94 3 Problem Formulation

95 Consider a robot \mathcal{R} observing a static scene composed of specific objects o_j sampled from some
 96 distribution of objects \mathcal{O} . The objects divide the workspace into occupied space \mathcal{W}_{occ} and free
 97 space $\mathcal{W}_{free} = \mathcal{W} \setminus \mathcal{W}_{occ}$. The robot has access to a training subset of \mathcal{O} beforehand, but does not
 98 know the specific objects o_j in the current scene.

99 The robot observes the scene with two distinct sensing modalities. In the visual modality, the robot
 100 views the scene from a stationary RGBD camera receiving color depth images Im . Due to sensor
 101 noise and occlusion these depth images offer an incomplete and noisy measurement on the full
 102 region of \mathcal{W}_{occ} occupied by the obstacles. From the camera image we assume the scene can be
 103 segmented into distinct objects from \mathcal{O} .

104 For the tactile modality, consider a robot that is able to sense *if* it has made contact with any object,
 105 but not *where* along the robot surface the contact was made. We assume the contact does not move
 106 the objects. For a configuration in configuration space $q \in \mathcal{C}$, let $\mathcal{R}(q) \subset \mathcal{W}$ denote the region of
 107 workspace occupied by the robot. A robot that has visited configurations $\{q_1, q_2, \dots\} = Q_{free} \subset \mathcal{C}$
 108 without observing contact can carve out regions of known free space:

$$\bigcup_{q \in Q_{free}} \mathcal{R}(q) = \mathcal{W}_{known-free} \subset \mathcal{W}_{free} \quad (1)$$

109 For each configuration $q_{contact} \in Q_{contact}$ where contact is observed, there must be at least one
 110 object point in collision with the robot (and not in known freespace).

$$\forall q_{contact} \in Q_{contact} \exists p_{contact} \in (\mathcal{R}(q_{contact}) \setminus \mathcal{W}_{known-free}) : p_{contact} \in \mathcal{W}_{occ} \quad (2)$$

111 Using existing nomenclature, each such region is called a Collision Hypothesis Set (CHS) [36].

112 Our objective is to model the conditional occupancy $p(\mathcal{W}_{occ} | \mathcal{O}, Im, Q_{free}, Q_{contact})$. Specifically,
 113 we desire a stochastic function $g(\mathcal{O}, Im, Q_{free}, Q_{contact})$ which generates sample \mathcal{W}_{occ} as similar
 114 as possible to the true conditional distribution. Since the true conditional distribution is unknown,
 115 in practice we seek to minimize the distance of drawn samples to the ground truth scene.

116 **4 Method**

117 Our approach is to use a particle filter storing a collection of latent shapes. We first segment the
 118 scene into distinct objects, then use an existing shape completion neural network to draw latent
 119 shape samples ψ_j for each object from $p(\psi_j|\mathcal{O}, Im)$, initializing the particle filter. Each particle
 120 can be decoded into the objects in a scene, thus the collection of particles represents the belief
 121 $p(\mathcal{W}_{occ}|\mathcal{O}, Im)$. We propose Constrained LAtent Shape Projection (CLASP) as the measurement
 122 update, projecting these samples onto the constraints imposed by Q_{free} and $Q_{contact}$. Our method
 123 is shown in Fig. 2, where the trapezoids f_{enc} and f_{dec} are the encoders and decoders of PSSNet [5].

124 **4.1 Initial Belief**

125 The RGBD camera images are passed to a segmentation algorithm, which yields distinct pixel re-
 126 gions in the image corresponding to different objects o_j . For each object o_j the corresponding
 127 portion of the depth image is converted first to a point cloud, then voxelgrids of known-occupied
 128 and known-free space centered around the visible object points with a transform T_j mapping the
 129 voxelgrid to the workspace coordinates.

130 For each object o_j we use the Plausible Shape Sampling Network (PSSNet) [5] f to generate possible
 131 shape completions. PSSNet is structured as a variational autoencoder. An encoder f_{enc} maps the
 132 known-free and known-occupied voxelgrids to a mean and variance in latent space. A latent vector
 133 ψ can be sampled and passed to the decoder f_{dec} , which outputs a probability of occupancy for each
 134 voxel. Thresholding (e.g. $p > 0.5$ for each voxel) yields a completed shape.

135 An object o_j that is representable by f can be stored compactly as ψ_j such that $f_{dec}(\psi_j) = o_j$. The
 136 transform T_j maps the completed shape into the workspace frame. A world is composed of static
 137 objects $\{o_1, o_2, \dots\} \in \mathcal{O}$. A particle ϕ stores a specific world as a sequence of latent-space vectors
 138 $\{\psi_1, \psi_2, \dots\}$. We sample worlds conditioned on only the depth-image observation by independently
 139 sampling latent vectors of objects. The initial belief is a set of particles $\{\phi_1, \phi_2, \dots\} \in \Phi$ generated
 140 from the information from the depth camera before any robot motion.

141 **4.2 Projecting a single object**

142 Sampling particles using only camera information may yield worlds that are inconsistent with the
 143 robot contact information. For example, PSSNet may predict objects that extend far into occluded
 144 space that intersect regions the robot has moved through. Alternatively, PSSNet may predict objects
 145 that do not extend into occluded space, and so the robot may observe contact with no object to
 146 explain the collision. Predicting shapes from vision and robot contact in a single pass would require
 147 a dataset specific to each robot and a specific set of motions.

148 To resolve these inconsistencies, sampled particles are projected onto the constraints in the latent
 149 space of the shape completion network, shown in Fig. 2. For sample i of object j , ψ_j^i induces a
 150 workspace occupancy. Our constraints lie in the workspace, but we wish to project the latent space
 151 vector. Therefore, the projection is accomplished by optimizing a loss via gradient updates on ψ_j^i
 152 while holding f_{dec} fixed, mirroring the process of training a neural network but optimizing the input
 153 instead of the network weights. Consider the unthresholded voxelgrid with values between 0 and 1
 154 produced by the decoder: $f_{dec}(\psi_j^i) = W_j^i$.

155 We optimize the loss: $\mathcal{L}_{all} = \mathcal{L}_{free} + \mathcal{L}_{occ} + \mathcal{L}_{prior}$

156 The first term \mathcal{L}_{free} penalizes all voxels predicted above a threshold δ that are known to be free.

$$\mathcal{L}_{free} = \sum_{x,y,z} \begin{cases} \max(W_j^i(x, y, z) - \delta, 0) & \mathcal{W}_{known-free}(x, y, z) = 1 \\ 0 & \text{otherwise} \end{cases} \quad (3)$$

157 The second term \mathcal{L}_{occ} penalizes unexplained contact. Each contact $q_{contact}$ must be caused by some
 158 object s.t. $\mathcal{R}(q_{contact}) \cap \mathcal{W}_{occ} \neq \emptyset$, however it is not obvious which object is responsible for the
 159 contact, or which voxel of the object contacted the robot. During optimization we consider a specific
 160 assignment of $Q_{contact}^j$ to object o_j . We define the assignment process in Section 4.3. Because a
 161 single occupied voxel is enough to explain a contact, at each iteration the loss is optimized based on

162 the maximum prediction of occupancy overlapping with the collision hypothesis set.

$$\mathcal{L}_{occ} = \sum_{q \in Q_{contact}^j} 1 - \max(\mathcal{R}(q) \cdot W_j^i) \quad (4)$$

163 The final term \mathcal{L}_{prior} penalizes deviation of ψ from the original distribution predicted by the encoder.
164 Without this constraint, ψ can deviate arbitrarily, losing all dependence on the depth image and
165 even leaving the training domain of f_{dec} . This would produce completions that no longer look like
166 objects. \mathcal{L}_{prior} is weighted by α to maintain a similar magnitude of gradients to \mathcal{L}_{free} and \mathcal{L}_{occ} .

$$\mathcal{L}_{prior} = -\alpha \log(P(\psi_j^i | f_{enc}(Im))) \quad (5)$$

167 Sampling the occupancy for a specific object o_j given Im , Q_{free} and $Q_{contact}^j$ is thus accomplished
168 by sampling a ψ_j and optimizing until the constraints are satisfied. Projection can fail if an iteration
169 limit, set to 100 steps, is reached without satisfying the constraints. For practical efficiency this
170 failure can sometimes be detected early when gradient updates no longer change the loss and the
171 constraints are not satisfied. We use Adam [37] for optimization with a learning rate of 0.01.

172 4.3 Multi-object completion

173 CLASP stores an assignment of each $q_{contact}$ to a particular object o_j for each full-scene particle
174 i . When a measurement contains a new contact $q_{contact}$, it is assigned to a specific object for each
175 sampled particle i as follows. First, the output of our shape completer is a finite-sized voxelgrid,
176 typically smaller than the full scene. The new $q_{contact}$ cannot be assigned to any object j where
177 $\mathcal{R}(q_{contact})$ lies entirely outside the output region of the decoder $f_{dec}(\psi_j)$. Next, for each remaining
178 j , a projection is attempted for each ψ_j^i to satisfy the new $q_{contact}$. If all attempts fail, we assume
179 this new $q_{contact}$ was not caused by object j . For each full-scene particle i , a specific assignment of
180 $q_{contact}$ is randomly and uniformly selected from the remaining objects j that could possibly explain
181 the contact.

182 5 Experiments

183 We evaluated scenarios of different objects to determine if CLASP improves the estimate of the
184 scene using robot contact information. We tested ablations of CLASP to evaluate the importance
185 of the latent prior and constraint satisfaction. We also tested alternative approaches to CLASP that
186 did not rely on projection. Finally, we compared CLASP on two different network architectures and
187 trained on multiple datasets. We trained separate instances of PSSNet [5] on Axis-Aligned Boxes
188 (AAB), YCB [38], and ShapeNet mugs [39] (training details in Section A.1).

189 5.1 Robot Contacts

190 **Simulation:** To generate contact measurements we moved the right arm of a robot composed of two
191 Kuka iiwa arms with Robotiq 3-finger grippers. We generated scenes by manually placing simulated
192 objects from AAB, YCB, or ShapeNet on a virtual table at about camera height. The known voxels
193 were passed to our trained PSSNet to generate a set of possible completions.

194 We generated robot motions to gather information by moving near and sometimes contacting the
195 objects using the procedure described in the appendix A.2. The first motion typically sweeps known
196 free space rather than making contact. The second or third motion intentionally makes contact with
197 the object.

198 **Live Robot:** The physical kinematics of our robot matched the simulated robot. A Kinect depth
199 camera mounted at the “head” position generated the RGBD images. A calibrated motion capture
200 system provided transforms between the Kinect and robot frames. We segmented the RGB image
201 using the CSAIL semantic-segmentation-pytorch library [40] which we retrained on YCB objects.
202 Each segmentation was converted into a voxelgrid and fed to PSSNet as in simulation to generate
203 sample worlds. The same procedure was used to generate robot motions as in simulation.

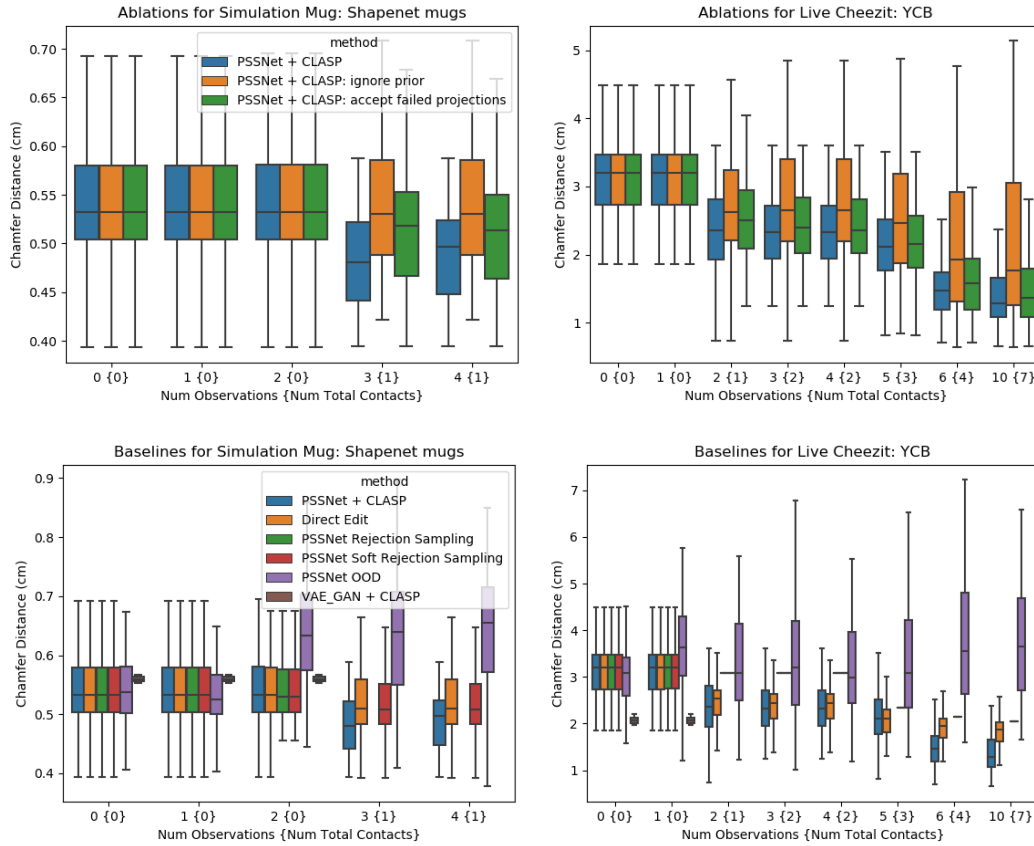


Figure 3: Boxplots showing the Chamfer Distance from sampled particles to ground truth. The mean, middle quartiles (boxed colored region), and outer quartiles excluding outliers are shown. Rejection Sampling and VAE_GAN occasionally produced no valid shapes, in which case no box is displayed.

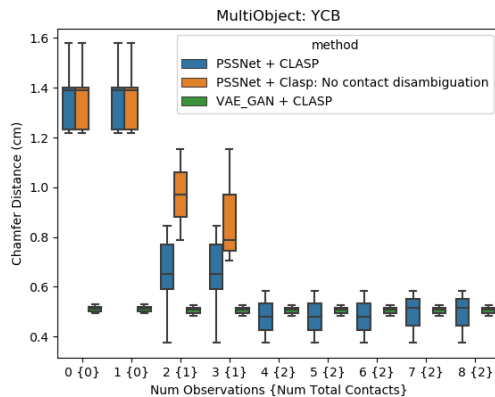


Figure 4: Boxplot results for the multiobject scene. PSSNET + CLASP: NO CONTACT DISAMBIGUATION fails to project any samples for observations 4 and beyond, so there is not corresponding box.

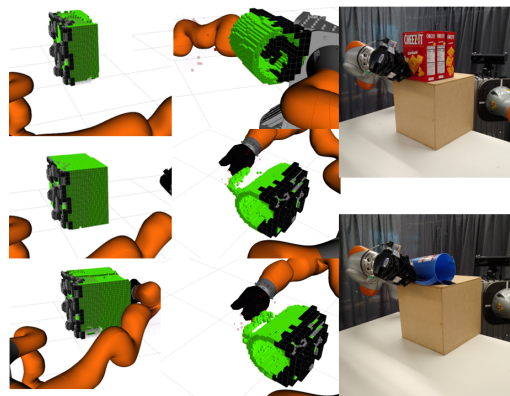


Figure 5: The Deep Cheezit (left), Mug (Middle), Live Cheezit and Live Pitcher (Right) scenes. The occupied (black) and known free (not shown) voxels from vision with contact (transparent red) and robot free (not shown) voxels are all used by CLASP to generate completed shapes (green).

204 On the live robot, contact was determined at each configuration by checking if the measured external
 205 torque exceeded a threshold of $2Nm$ per joint. This threshold was large enough to avoid generating
 206 false positives while remaining sensitive enough to detect contact with secured objects.

207 5.2 Scenes

208 We tested four scenes in simulation and two on
 209 the live robot. Each scene consisted of a single
 210 object secured to the table in front of the robot.
 211 We also tested a scene of multiple YCB objects
 212 (Fig. 6). Contacts occurred with occluded sec-
 213 tions of the objects, with examples shown in
 214 Fig. 5. In both simulation and the live robot,
 215 table occupancy was not considered when evalu-
 216 ating the quality of the completions.

217 **Simulated Scenes:** The first pair of scenarios
 218 used a single YCB Cheezit box (Shallow) and
 219 a stack of three Cheezit boxes (Deep). These
 220 setups generated similar depth images but dif-
 221 ferent ground truth shapes. Both used net-
 222 works trained on the AAB dataset. The Simu-
 223 lated Pitcher from YCB was positioned with the
 224 handle occluded from view and used networks
 225 trained on the full YCB dataset. The Simulated
 226 Mug from ShapeNet also had the handle oc-
 227 cluded from view and used networks trained on all mugs in ShapeNet. The handles on these objects
 228 were localized through contact.

229 **Live Scenes:** The Live Cheezit also consisted of a stack of three boxes, and again the Live YCB
 230 Pitcher had the handle occluded. Both scenes used networks trained on the full YCB dataset. The
 231 Cheezit boxes were attached together and the pitcher was taped to prevent motion during contact.
 232 Simulated objects were manually aligned to the live scene to approximate the ground truth of live
 233 objects and were used for evaluation.

234 5.3 Baselines

235 We compared our proposed method to the following alternatives. **DIRECT EDIT** directly adds or
 236 removes voxels to satisfy the contact information. **REJECTION SAMPLING** samples latent space
 237 vectors from the distribution predicted by the encoder, then decoded these into 3D objects and
 238 rejected any samples not satisfying the contact or free space constraints. **SOFT REJECTION SAM-**
 239 **PLING** selects and then Directly Edits the least violating samples in the cases where all samples are
 240 rejected. For **OOD (Direct Out-Of-Distribution Prediction)**, we combined the visual known free and
 241 occupied voxels with the contact known free and occupied voxels as input to PSSNet. While other
 242 methods did not have access to the true contact point, we allowed this method this advantage and
 243 added the true contact point directly to the known occupied voxels from the depth image. Finally,
 244 while other approaches used the PSSNet shape completion network, we tested using a VAE.GAN
 245 [17] network. This network tends to produce better average but less diverse samples. Section A.3
 246 describes these baselines in more detail.

247 We also tested ablations of our method. **CLASP: IGNORE PRIOR** tested removing the loss term
 248 \mathcal{L}_{prior} . **CLASP: ACCEPT FAILED PROJECTIONS** tested accepting all projections, even those that
 249 do not satisfy the contact constraints. To test our contact assignment in the multi-object case (Sec-
 250 tion 4.3), **CLASP: NO CONTACT DISAMBIGUATION** determined if a projection of latent ψ_j could
 251 satisfy each new $q_{contact}$ as in CLASP, then assigned each $q_{contact}$ to all feasible objects j . This
 252 resulted in scenes explaining a single $q_{contact}$ with multiple objects.

253 100 particles were sampled in each method, with the threshold of \mathcal{L}_{free} set at $\delta = 0.4$ and the
 254 weighting of \mathcal{L}_{prior} set at $\alpha = 0.01$.

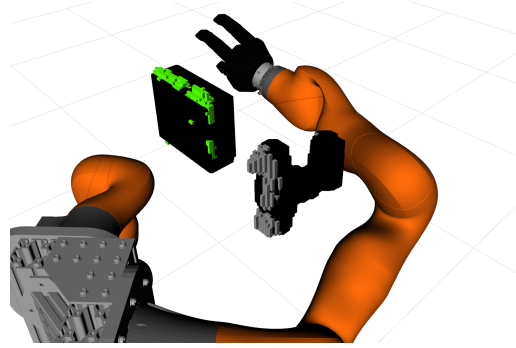


Figure 6: The multiobject scene with the YCB Cheezit box and Drill on a table (not shown) just before the first contact.

255 5.4 Results

256 Single-object scenes were tested on all baselines and ablations using PSSNet trained with the ap-
257 propriate dataset. Fig. 3 compares the Chamfer Distance (CD) [41] of each accepted sample to
258 the ground truth for selected scenes. Plots for all scenes are shown in Section A.3. We consider a
259 different analysis in Section A.4. We find that REJECTION SAMPLING often performs well during
260 the first few observations with zero or one contact. However, REJECTION SAMPLING soon fails
261 to return any valid samples with two or more contacts. We see that OOD produces completions
262 that are typically much worse than the original completions from only vision. The initial estimate
263 (observation 0) from VAE_GAN are hit-or-miss. All networks saw the YCB Pitcher during training,
264 and VAE_GAN recalls the pitcher more accurately than PSSNet during testing to the point where
265 contacts are unnecessary. However, the recall of VAE_GAN in the other ambiguous scenarios is
266 worse than PSSNet and projection to the contact constraints often fails, leaving no sampled shapes.
267 The results justify our choice of PSSNet over VAE_GAN for CLASP, as VAE_GAN is unable to
268 sufficiently adjust to the contact information.

269 Considering ablations of CLASP, ACCEPT FAILED PROJECTIONS performed as well initially (when
270 no projections fail) and significantly worse as the number of observations increases. Ignoring the la-
271 tent prior during projection also performs worse and occasionally produces shapes that qualitatively
272 look less like objects compared to completions from other methods.

273 Across all scenes, CLASP performed similarly to the best of all other methods with 0 or 1 contacts
274 and the best with multiple contacts. CLASP successfully used the robot contact information in all
275 scenarios to reduce the CD between the predicted and ground truth shapes in all scenarios. Robot
276 measurements with a contact typically caused a larger reduction in CD than measurements with only
277 freespace information. Numerically, the CD reduced most in the Cheezit scenarios, with a reduction
278 of the mean from $0.5cm$ to $0.1cm$ for the Shallow and from $0.14cm$ to $0.08cm$ for the Deep. The
279 CD reduction in the pitcher and mug scenarios was significantly smaller, as the general shape of
280 the pitcher and mug could be predicted from the image. The prediction of the occluded handle
281 was improved with contact. The trend of improvement in the live scenes matched the simulation.
282 However, the numeric error of the live scenes was much larger, perhaps caused by imperfect transfer
283 of the learned shape completer from training in simulation to prediction on live Kinect data as well
284 as imperfect alignment of the robot frame to the Kinect frame.

285 In the multi-object scene (Fig. 4) the VAE_GAN method achieves a better completion from the
286 RGBD data, but our proposed method produces better samples after 2 contacts. Our proposed con-
287 tact assignment (Section 4.3) outperforms naively satisfying the contacts whenever possible.

288 6 Discussion and Conclusion

289 While we model CLASP using a particle filter and would like to have the Bayesian estimate of the
290 scene given all observations, we acknowledge many non-Bayesian approximations. Particle filters
291 approximate Bayes filters, but the 100 particles we sample may not be a sufficient coverage of the
292 latent shape space. CLASP projects samples, which does not preserve Bayesian estimates.

293 While our shape network uses voxelgrids, implicit representations have recently become popular
294 and offer advantages worth considering. Currently shape completion networks produce the most
295 visually pleasing results when trained on a single object, visually decent results when trained on a
296 single class of objects, and poor results when trained on large diverse shape datasets. In order to
297 be practically applicable to robots, shape completion must handle a wide variety of objects. Shape
298 completion is rarely the end goal, but rather a tool robots can use to aid in tasks. Choices of correct
299 metrics and refinements to CLASP ultimately depend on the specific downstream application.

300 We demonstrated a method for estimating shape completions initialized with purely RGBD visual
301 data, then updated from observations of a robot arm moving through unknown regions and sensing
302 contact. We stored the belief of the scene as a particle filter of latent vectors from a shape com-
303 pletion network and used CLASP to enforce shape consistency with the robot observations. Most
304 importantly, we showed that CLASP improves the estimate of object shape using these contact ob-
305 servations. Our results further showed that CLASP performs better than ablations of CLASP and
306 alternative methods. We hope CLASP will be used within a larger robotics framework where rea-
307 soning over environment uncertainty based on shape priors aids in accomplishing larger goals.

References

- 308
309 [1] M. Klingensmith, M. Koval, S. Srinivasa, N. Pollard, and M. Kaess. The manifold particle filter for state
310 estimation on high-dimensional implicit manifolds, April 2016.
- 311 [2] K. Desingh, S. Lu, A. Pipari, and O. C. Jenkins. Factored pose estimation of articulated objects using
312 efficient nonparametric belief propagation. In *ICRA*, 2019.
- 313 [3] M. A. Lee, Y. Zhu, P. Zachares, M. Tan, K. Srinivasan, S. Savarese, L. Fei-Fei, A. Garg, and J. Bohg.
314 Making sense of vision and touch: Learning multimodal representations for contact-rich tasks. *IEEE*
315 *Transactions on Robotics*, 2020. doi:10.1109/TRO.2019.2959445.
- 316 [4] M. C. Koval, N. S. Pollard, and S. S. Srinivasa. Pose estimation for planar contact manipulation with
317 manifold particle filters. *IJRR*, 2015.
- 318 [5] B. Saund and D. Berenson. Diverse plausible shape completions from ambiguous depth images. *CoRL*,
319 2020.
- 320 [6] Zhirong Wu, S. Song, A. Khosla, Fisher Yu, Linguang Zhang, Xiaoou Tang, and J. Xiao. 3D shapenets:
321 A deep representation for volumetric shapes. In *CVPR*, 2015.
- 322 [7] C. B. Choy, D. Xu, J. Gwak, K. Chen, and S. Savarese. 3D-R2N2: A unified approach for single and
323 multi-view 3D object reconstruction. In *ECCV*, 2016.
- 324 [8] R. Girdhar, D. Fouhey, M. Rodriguez, and A. Gupta. Learning a predictable and generative vector repre-
325 sentation for objects. In *ECCV*, 2016.
- 326 [9] J. Wu, C. Zhang, X. Zhang, Z. Zhang, W. T. Freeman, and J. B. Tenenbaum. Learning shape priors for
327 single-view 3D completion and reconstruction. In *ECCV*, 2018.
- 328 [10] J. Wu, Y. Wang, T. Xue, X. Sun, W. T. Freeman, and J. B. Tenenbaum. MarrNet: 3D Shape Reconstruction
329 via 2.5D Sketches. In *Advances In Neural Information Processing Systems*, 2017.
- 330 [11] M. Michalkiewicz, E. Belilovsky, M. Baktashmotlagh, and A. Eriksson. A simple and scalable shape
331 representation for 3D reconstruction. *arXiv*, 2020.
- 332 [12] C. Wen, Y. Zhang, Z. Li, and Y. Fu. Pixel2mesh++: Multi-view 3D mesh generation via deformation.
333 *ICCV*, 2019.
- 334 [13] H. Xie, H. Yao, X. Sun, S. Zhou, and S. Zhang. Pix2vox: Context-aware 3d reconstruction from single
335 and multi-view images. In *ICCV*, 2019.
- 336 [14] H. Fan, H. Su, and L. J. Guibas. A point set generation network for 3D object reconstruction from a single
337 image. *CVPR*, 2017.
- 338 [15] Y. Yu, Z. Huang, F. Li, H. Zhang, and X. Le. Point encoder gan: A deep learning model for 3D point
339 cloud inpainting. *Neurocomputing*, 2020.
- 340 [16] B. Yang, S. Rosa, A. Markham, N. Trigoni, and H. Wen. Dense 3D object reconstruction from a single
341 depth view. In *TPAMI*, 2018.
- 342 [17] J. Wu, C. Zhang, T. Xue, B. Freeman, and J. Tenenbaum. Learning a probabilistic latent space of object
343 shapes via 3D generative-adversarial modeling. In *NIPS*, 2016.
- 344 [18] L. Roldão, R. de Charette, and A. Verroust-Blondet. 3d semantic scene completion: a survey. *ArXiv*,
345 abs/2103.07466, 2021.
- 346 [19] J. Jung. Active shape completion using tactile glances. In *Master's Thesis: Lulea University of Technol-*
347 *ogy*, 2019.
- 348 [20] P. J. Besl and N. D. McKay. A method for registration of 3-D shapes. *TPAMI*, 1992.
- 349 [21] H. Yang, J. Shi, and L. Carlone. TEASER: Fast and Certifiable Point Cloud Registration. *arXiv*, 2020.
- 350 [22] V. Narayanan and M. Likhachev. Discriminatively-guided deliberative perception for pose estimation of
351 multiple 3D object instances. In *RSS*, June 2016.
- 352 [23] X. Deng, A. Mousavian, Y. Xiang, F. Xia, T. Bretl, and D. Fox. PoseRBPF: A rao-blackwellized particle
353 filter for 6D object pose estimation. In *RSS*, 2019.

- 354 [24] T. Hodan, D. Barath, and J. Matas. EPOS: Estimating 6d pose of objects with symmetries. In *IEEE/CVF*
355 *Conference on Computer Vision and Pattern Recognition*, 2020.
- 356 [25] W. Yuan, S. Dong, and E. Adelson. Gelsight: High-resolution robot tactile sensors for estimating geome-
357 try and force. *Sensors*, 2017.
- 358 [26] A. Alspach, K. Hashimoto, N. Kuppaswamy, and R. Tedrake. Soft-bubble: A highly compliant dense
359 geometry tactile sensor for robot manipulation. In *RoboSoft*, 2019.
- 360 [27] S. Luo, W. Mou, K. Althoefer, and H. Liu. Localizing the object contact through matching tactile features
361 with visual map. *ICRA*, 2015.
- 362 [28] W. Yuan, Y. Mo, S. Wang, and E. Adelson. Active clothing material perception using tactile sensing and
363 deep learning. *ICRA*, 2018.
- 364 [29] N. Kuppaswamy, A. Alspach, A. Uttamchandani, S. Creasey, T. Ikeda, and R. Tedrake. Soft-bubble
365 grippers for robust and perceptive manipulation. In *IROS*, 2020.
- 366 [30] Y. Li, J.-Y. Zhu, R. Tedrake, and A. Torralba. Connecting touch and vision via cross-modal prediction. In
367 *CVPR*, 2019.
- 368 [31] W. Agnew, C. Xie, A. Walsman, O. Murad, C. Wang, P. Domingos, and S. S. Srinivasa. Amodal 3d
369 reconstruction for robotic manipulation via stability and connectivity. *CoRL*, 2020.
- 370 [32] Q. Lu, M. V. der Merwe, and T. Hermans. Multi-fingered active grasp learning. *IROS*, 2020.
- 371 [33] Z. Yi, R. Calandra, F. Veiga, H. van Hoof, T. Hermans, Y. Zhang, and J. Peters. Active tactile object
372 exploration with gaussian processes. *IROS*, 2020.
- 373 [34] S. Dragiev, M. Toussaint, and M. Gienger. Uncertainty aware grasping and tactile exploration. *ICRA*,
374 2013.
- 375 [35] S. Wang, J. Wu, X. Sun, W. Yuan, W. Freeman, J. Tenenbaum, and E. Adelson. 3d shape perception from
376 monocular vision, touch, and shape priors. In *IROS*, 2018.
- 377 [36] B. Saund and D. Berenson. Motion planning for manipulators in unknown environments with contact
378 sensing uncertainty. In *ISER*, 2018.
- 379 [37] D. Kingma and J. Ba. Adam: A method for stochastic optimization. *International Conference on Learning*
380 *Representations*, 2014.
- 381 [38] B. Calli, A. Singh, J. Bruce, A. Walsman, K. Konolige, S. Srinivasa, P. Abbeel, and A. M. Dollar. Yale-
382 CMU-Berkeley dataset for robotic manipulation research. *IJRR*, April 2017.
- 383 [39] A. X. Chang, T. Funkhouser, L. Guibas, P. Hanrahan, Q. Huang, Z. Li, S. Savarese, M. Savva, S. Song,
384 H. Su, J. Xiao, L. Yi, and F. Yu. ShapeNet: An Information-Rich 3D Model Repository. Technical Report
385 arXiv, 2015.
- 386 [40] B. Zhou, H. Zhao, X. Puig, T. Xiao, S. Fidler, A. Barriuso, and A. Torralba. Semantic understanding of
387 scenes through the ade20k dataset. *International Journal on Computer Vision*, 2018.
- 388 [41] H. Barrow, J. Tenenbaum, R. Bolles, and H. Wolf. Parametric correspondence and chamfer matching:
389 Two new techniques for image matching. *IJCAI*, 1977.
- 390 [42] F. S. Nooruddin and G. Turk. Simplification and repair of polygonal models using volumetric techniques.
391 *IEEE Transactions on Visualization and Computer Graphics*, 2003.
- 392 [43] P. Min. binvox. <http://www.patrickmin.com/binvox>, 2004 - 2020.
- 393 [44] D. Berenson, S. Srinivasa, and J. Kuffner. Task space regions: A framework for pose-constrained manip-
394 ulation planning. *International Journal of Robotics Research*, 30(12):1435 – 1460, October 2011.
- 395 [45] B. Saund, S. Chen, and R. Simmons. Touch based localization of parts for high precision manufacturing.
396 In *ICRA*, 2017.
- 397 [46] A. Hermann, F. Drews, J. Bauer, S. Klemm, A. Roennau, and R. Dillmann. Unified GPU voxel collision
398 detection for mobile manipulation planning. In *IROS*, 2014.

APPENDIX

A Experiment Details

A.1 Shape Network Training

We generated 3 distinct datasets of voxelized objects with size 64^3 from random axis-aligned boxes (AAB), YCB objects [38], and ShapeNet mugs [39]. Boxes for AAB had width, depth, and height uniformly sampled with 2 to 41 voxels. For YCB and ShapeNet we generated ground truth voxelgrids centered on the object with different rotations using `binvox` [42, 43]. For YCB we applied all 15 degree increment rotations about both the vertical and a horizontal axis. For Shapenet we applied all 5 degree increment rotations about the vertical axis.

During each epoch of training, each voxelgrid was augmented with translations sampled uniformly from -10 to 10 voxels in each direction. The 2.5D “known occupied” and “known free” voxelgrids were generated assuming a sensor looking down the x-direction. Sensor noise was simulated by sampling IDD 0-mean 2cm-std. deviation gaussian random noise in a depth image of 16x16, scaling that depth image to 64x64 using bilinear interpolation, then applying that noise to the x-direction of the known occupied and free voxelgrids.

We trained separate instances of PSSNet [5] on AAB, YCB, and Shapenet mugs, training for at least 100 epochs (~ 1 day), and used the iteration with minimal loss for experiments.

A.2 Robot Motion Generation

The following procedure was used to generate the robot motion which in turn generated the contact and freespace observations. The robot began each trial with a roadmap: a graph of nodes corresponding to configurations, and edges of robot motions connecting the nodes. Each scene contained a Goal Generator function, which mapped the completed objects to a goal Task-Space Region (TSR) [44]. Using 10 sampled worlds, there were a corresponding 10 separate TSRs. In the scene above, the Goal Generator took the mean of the completed object points, and generated a TSR centered 10cm back in the occluded region.

At each iteration if the robot did not currently satisfy any TSR, approximately 80 configurations were sampled from each TSR and added to the roadmap, and the robot would attempt to traverse the roadmap to the closest configuration in a TSR. If the robot satisfied all TSRs the task was considered complete.

If instead the robot satisfied at least one but not all TSRs, the robot took an information gathering action. For each outgoing edge of the robot’s node on the roadmap, the Information Gain (IG) was calculated from the existing particles using the method in [45]. The robot took the action with highest information gain, which often (intentionally) contacted an object. The belief was updated and the next iteration began.

Detecting contact in simulation: The voxelized robot was computed using `GpuVoxels` [46] with a much larger 256^3 voxelgrid with 1cm voxel side lengths. This robot voxelgrid was converted to an occupied point cloud, then transformed to the object frame, and converted into a voxelgrid matching the size and position of the depth image voxelgrid. Contact was determined by checking for overlap between the robot and object voxelgrid. For each configuration visited not in contact, the voxelized robot was added to the known freespace. Each configuration in contact generated a Collision Hypothesis Set, added to $Q_{contact}$.

A.3 Baselines and All Results

Here we describe several baselines in more details. Figures 7 and 8 show all ablations and baselines for all scenes described in Section 5.

DIRECT EDIT: Perhaps the simplest method, in this baseline we apply the contact information directly to the voxelgrid. We sample latent vectors directly from the visual prior and decode into voxelgrid shapes as in other methods. At each measurement the known-free voxels from contact information are directly removed from predicted voxelgrids. In our problem formulation the true

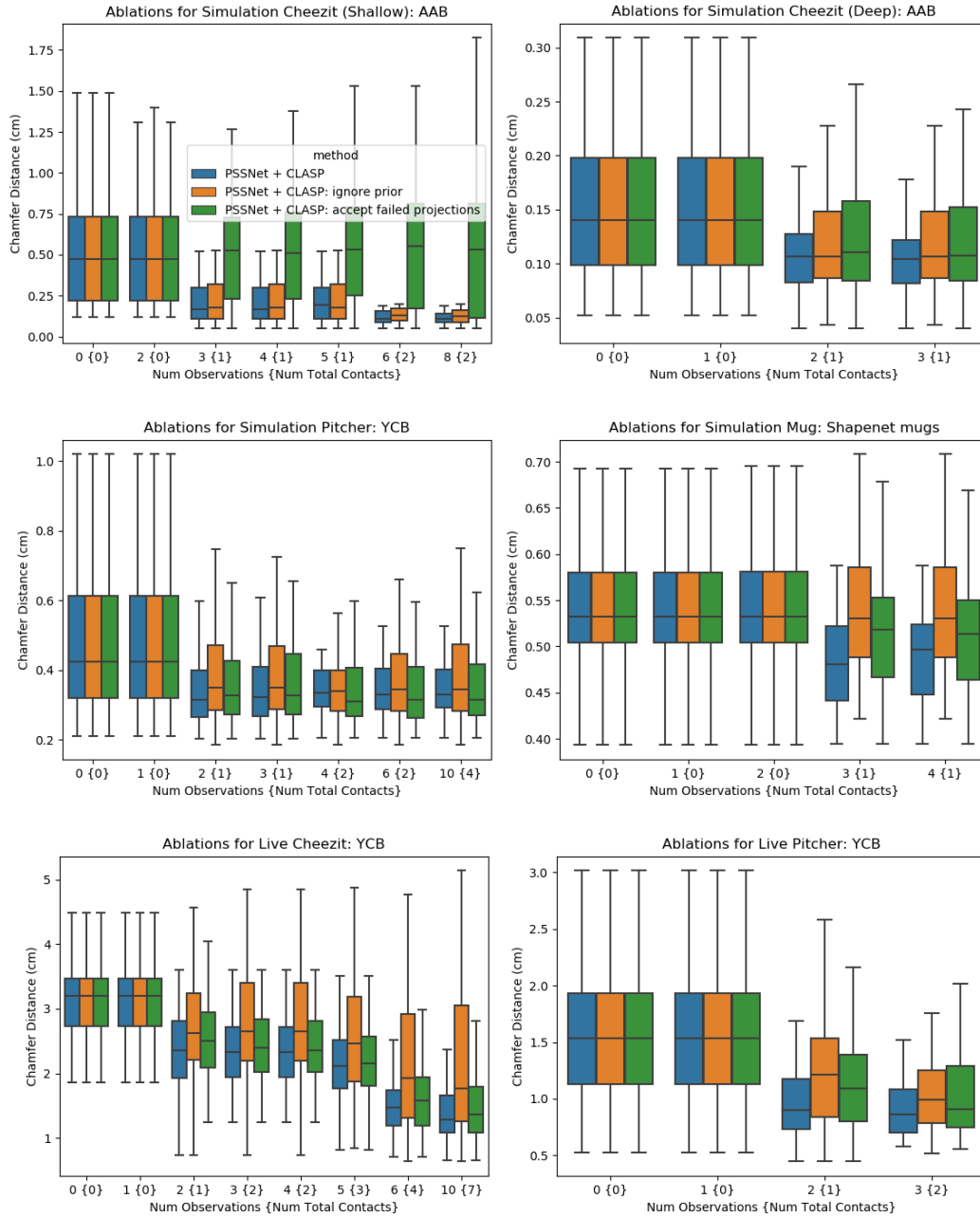


Figure 7: Ablation boxplots showing the Chamfer Distance from sampled particles to ground truth. The mean, middle quartiles (boxed colored region), and outer quartiles excluding outliers are shown.

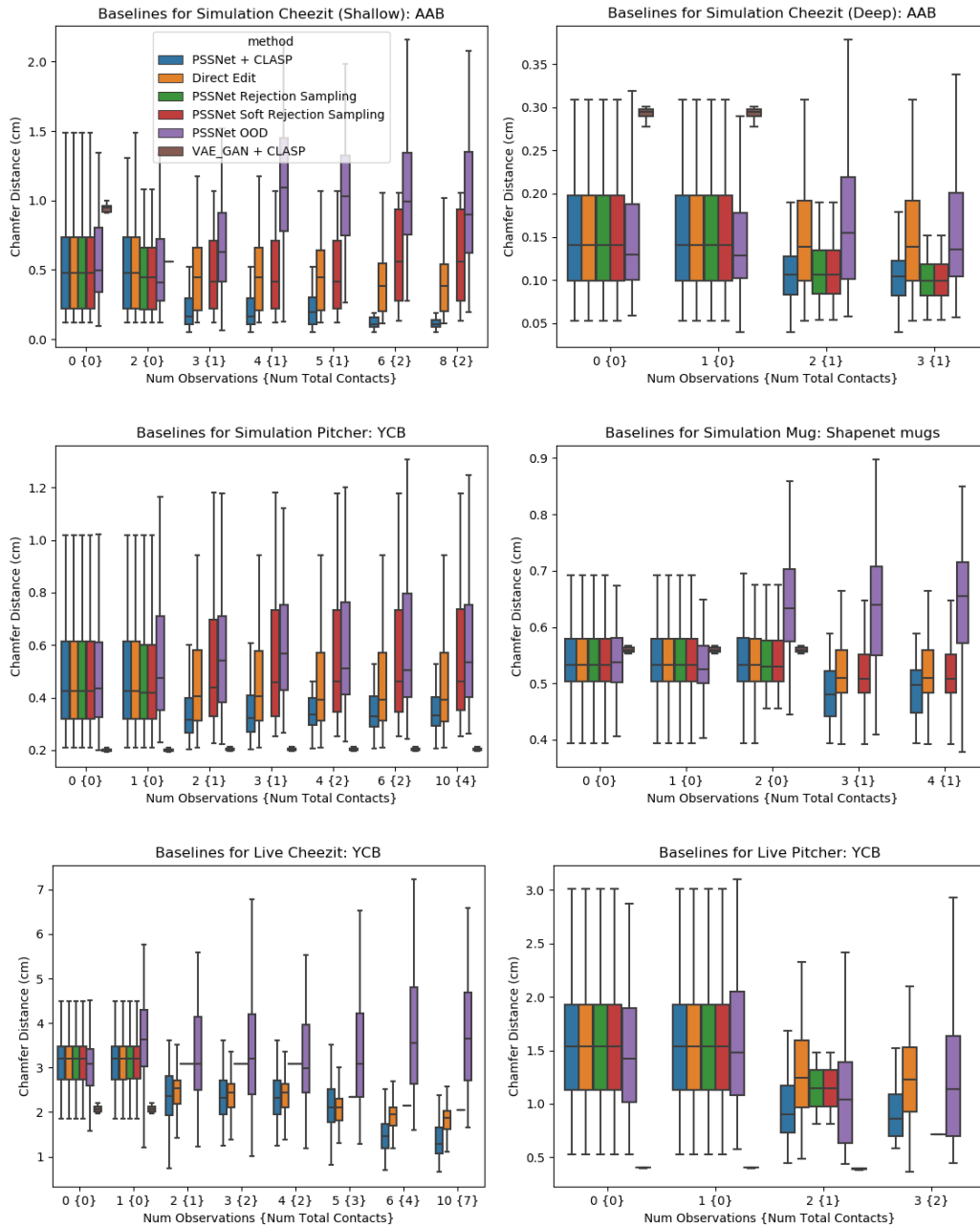


Figure 8: Baseline boxplots showing the Chamfer Distance from sampled particles to ground truth. The mean, middle quartiles (boxed colored region), and outer quartiles excluding outliers are shown. Rejection Sampling and VAE_GAN occasionally produced no valid shapes, in which case no box is displayed.

447 contact voxels are not known, however for this baseline we break that assumption and use the true
 448 object voxels that came in contact with the robot. These contacted voxels are added directly to the
 449 predicted shapes.

450 **SOFT REJECTION SAMPLING:** Rejection Sampling samples latent vectors from the distribution
 451 predicted by the encoder using purely visual data. To condition on the contact information, Rejec-
 452 tion Sampling simply discards all samples that do not satisfy the contact constraints. To overcome
 453 the limitation that after 1 or two contacts no samples are accepted, reviewers suggested this softer
 454 implementation. In the SOFT REJECTION SAMPLING baseline we perform rejection sampling, sav-
 455 ing rejected samples but marking them as invalid. If all samples are marked as invalid, we select
 456 all samples that violate the fewest number of constraints. This count includes each CHS without an
 457 occupied voxel, and each known-free voxel which the sample occupies. We then apply the DIRECT
 458 EDIT method on each selected sample to enforce that each sample satisfies all known constraints.

459 **OOD (Direct Out-Of-Distribution Prediction):** Our neural network accepts known-free and
 460 known-occupied voxels as input. During training the known-free and known-occupied solely came
 461 from vision, however contact information provides the same type of known-occupied and known-
 462 free information. This baseline takes the union of the known-free voxels from vision and robot
 463 motion as the known-free input, and uses the union of known-occupied voxels from vision and con-
 464 tact as the known-occupied input. In our problem formulation the true contact voxels are not known,
 465 however for this baseline (as in DIRECT EDIT) we break that assumption and use the true object
 466 voxels that came in contact with the robot. This baseline was suggested by many colleagues in early
 467 discussions of the paper. It is perhaps not surprising that this baseline performs poorly, as the com-
 468 bined information from vision and contact is out of distribution from all data on which the network
 469 was trained.

470 A.4 Likelihood Results

471 We consider an alternative analysis of the experiment data presented in Section 5.4. Given that we
 472 model scenes by sampling shapes in a particle filter, we consider the likelihood of the ground truth
 473 scene given the particles. Since a particle filter models discrete samples, none of which will exactly
 474 match the ground truth, we apply a kernel to our particles in workspace. Specifically, we apply a
 475 non-normalized kernel based on the Chamfer Distance between two shapes s_1, s_2 :

$$k(s_1, s_2) = \frac{1}{CD(s_1, s_2)} \quad (6)$$

476 The (non-normalized) likelihood of a particular scene occupancy s under the belief of n particles Φ
 477 is then

$$p(s|\Phi) = \sum_{\phi \in \Phi} \frac{1}{n} k(f_{dec}(\phi), s) \quad (7)$$

478 where $f_{dec}(\phi)$ decodes all latent shape vectors $\psi \in \phi$ into a scene.

479 We plot the likelihood of the true scene in Fig. 9 and Fig. 10, and find similar trends as in Section 5.4.
 480 The magnitude of the likelihood is not meaningful, however the relatively likelihoods between the
 481 methods are. Initially methods perform similarly, except VAE_GAN which is either better or worse
 482 than other methods. With contact and freespace observations, our proposed CLASP with PSSNet
 483 tends to increase the likelihood of the ground truth scene, while VAE_GAN tends to decrease the
 484 likelihood.

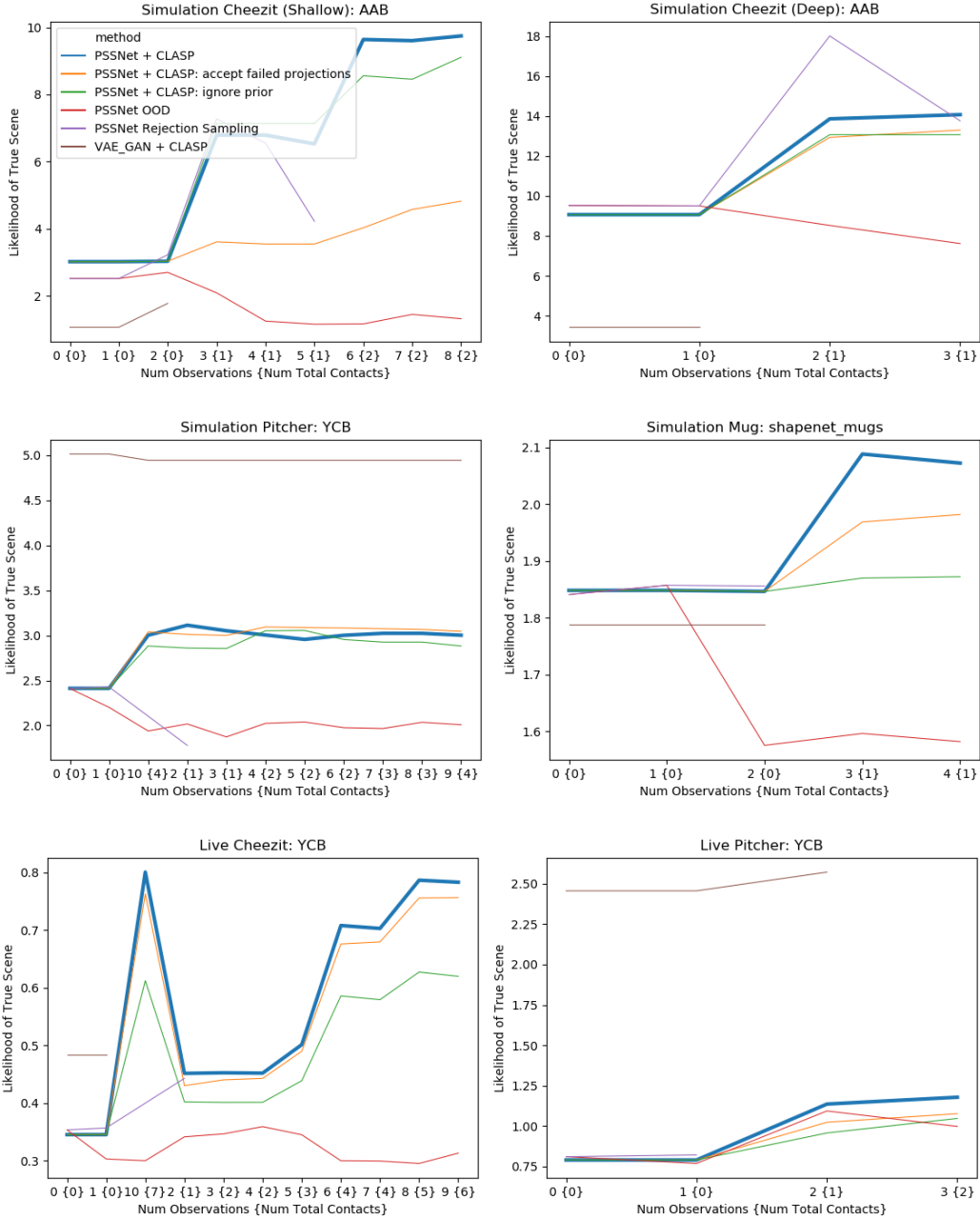


Figure 9: Plots of likelihoods of CLASP and baselines under the particle filter belief and kernel function.

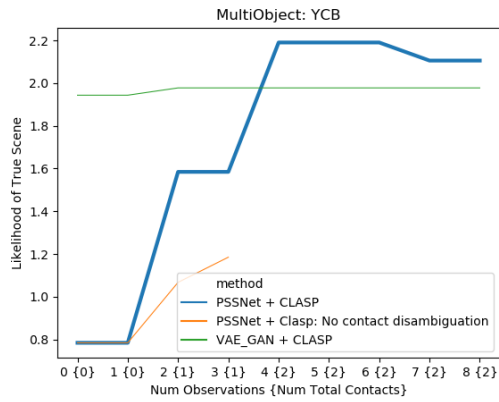


Figure 10: Likelihood of CLASP and baselines for the multiobject scene