

404 A Method Details

405 We included the code for both our simulated and real-world experiments for reference. Please find
 406 it in the supplementary material under `iw_code`. Algorithms 1 and 2 describe our warp learning and
 407 inference.

Algorithm 1 Warp Learning

Input: Meshes of K example object instances $\{\text{obj}_1, \text{obj}_2, \dots, \text{obj}_K\}$.
Output: Canonical point cloud, vertices and faces and a latent space of warps.
Parameters: Smoothness of CPD warping α and number of PCA components L .

- 1: $\text{PCD} = \langle \text{SampleS}(\text{obj}_i) \rangle_{i=1}^K$. ▷ Sample a small point cloud per object (Appendix A.1).
- 2: $C = \text{SelectCanonical}(\text{PCD})$. ▷ Select a canonical object with index C (Appendix A.2).
- 3: $\text{canon} = \text{Concat}(\text{obj}_C.\text{vertices}, \text{SampleL}(\text{obj}_C))$. ▷ Use both vertices and surface samples.
- 4: **for** $i \in \{1, 2, \dots, K\}, i \neq C$ **do**
- 5: $W_{C \rightarrow i} = \text{CPD}(\text{canon}, \text{PCD}_i, \alpha)$. ▷ Coherent Point Drift warping (Section 3).
- 6: **end for**
- 7: $D_W = \{\text{Flatten}(W_{C \rightarrow i})\}_{i=1, i \neq C}^K$. ▷ Dataset of displacements of canon.
- 8: $\text{PCA} = \text{FitPCA}(D_W, \text{n.components} = L)$. ▷ Learn a latent space of canonical object warps.
- 9: **return** $\text{Canon}(\text{points} = \text{canon}, \text{vertices} = \text{obj}_C.\text{vertices}, \text{faces} = \text{obj}_C.\text{faces}), \text{PCA}$.

Algorithm 2 Warp Inference and Mesh Reconstruction

Input: Observed point cloud pcd , canonical object canon and latent space PCA .
Output: Predicted latent shape v and pose T .
Parameters: Number of random starts S , number of gradient descent steps T , learning rate η and object size regularization β .

- 1: $t_g = \frac{1}{|\text{pcd}|} \sum_{i=1}^{|\text{pcd}|} \text{pcd}_i$.
- 2: $\text{pcd} = \text{pcd} - t_g$. ▷ Center the point cloud.
- 3: **for** $i = 1$ **to** S **do**
- 4: $R_{\text{init}} = \text{Random initial 3D rotation matrix}$.
- 5: Initialize $v = (0 \ 0 \ \dots \ 0)$, $s = (1 \ 1 \ 1)$, $t_l = (0 \ 0 \ 0)$, $\hat{R} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \end{pmatrix}$.
- 6: Initialize Adam [36] with parameters v, s, t_l, r and learning rate η .
- 7: **for** $j = 1$ **to** T **do**
- 8: $\delta = \text{Reshape}(Wv)$.
- 9: $X = \text{canon.points} + \delta$. ▷ Warped canonical point cloud.
- 10: $R = \text{GramSchmidt}(\hat{R})$.
- 11: $X = (X \odot s) R_{\text{init}}^T R^T + t_l$. ▷ Scaled, rotated and translated point cloud.
- 12: $\mathcal{L} = \frac{1}{|\text{pcd}|} \sum_k^{|\text{pcd}|} \min_l^{|X|} \|\text{pcd}_k - X_l\|_2^2$. ▷ One-sided Chamfer distance.
- 13: $\mathcal{L} = \mathcal{L} + \beta \max_l^{|X|} \|X_l\|_2^2$. ▷ Object size regularization.
- 14: Take a gradient descent step to minimize \mathcal{L} using Adam.
- 15: **end for**
- 16: **end for**
- 17: Find parameters $v^*, s^*, t_l^*, R_{\text{init}}^*, R^*$ with the lowest final loss across $i \in \{1, 2, \dots, S\}$.
- 18: $X = \text{canon.points} + \text{Reshape}(Wv^*)$.
- 19: $X = (X \odot s^*)(R_{\text{init}}^*)^T (R^*)^T + t_l^* + t_g$. ▷ Complete point cloud in workspace coordinates.
- 20: $\text{vertices} = \langle X_1, X_2, \dots, X_{|\text{canon.vertices}|} \rangle$. ▷ First $|\text{canon.vertices}|$ points of X are vertices.
- 21: **return** $\text{Mesh}(\text{vertices} = \text{vertices}, \text{faces} = \text{canon.faces})$. ▷ Warped mesh.

408 A.1 Point Cloud Sampling

409 We use `trimesh`¹ to sample the surface of object meshes. The function
 410 `trimesh.sample.sample_surface_even` samples a specified number of points and then
 411 rejects points that are too close together. We sample 2k points for small point clouds (SampleS)
 412 and 10k point for large point clouds (SampleL).

413 A.2 Canonical Object Selection

414 Among the K example objects, we would like to find the one that is the easiest to warp to the other
 415 objects. For example, if we have ten examples of mugs, but only one mug has a square handle,
 416 we should not choose it as it might be difficult to warp it to conform to the round handles of the
 417 other nine mugs. We use Algorithm 3, which computes $K * K - 1$ warps and picks the object that
 418 warps to the other $K - 1$ objects with the lowest Chamfer distance. We also note an alternative and
 419 computationally cheaper algorithm from Thompson et al. [8], Algorithm 4. This algorithm simply
 420 finds the object that is the most similar to the other $K - 1$ objects without any warping.

Algorithm 3 Exhaustive Canonical Object Selection

Input: Point clouds of K training objects $\langle X^{(1)}, X^{(2)}, \dots, X^{(K)} \rangle$.
Output: Index of the canonical object.

```

1: for  $i = 1$  to  $K$  do
2:   for  $j = 1$  to  $K, j \neq i$  do
3:      $W_{i \rightarrow j} = \text{CPD}(X^{(i)}, X^{(j)})$  ▷ Warp point cloud  $i$  to point cloud  $j$ .
4:      $C_{i,j} = \frac{1}{|X^{(j)}|} \sum_{k=1}^{|X^{(j)}|} \min_{l=1}^{|X^{(i)}|} \|X_k^{(j)} - (X^{(i)} + W_{i \rightarrow j})_l\|_2^2$ 
5:   end for
6: end for
7: for  $i = 1$  to  $K$  do
8:    $C_i = \sum_{j=1, j \neq i}^K C_{i,j}$  ▷ Cumulative cost of point cloud  $i$  warps.
9: end for
10: return  $\arg \min_{i=1}^K C_i$  ▷ Pick point cloud that is the easiest to warp.

```

Algorithm 4 Approximate Canonical Object Selection [8]

Input: Point clouds of K training objects $\langle X^{(1)}, X^{(2)}, \dots, X^{(K)} \rangle$.
Output: Index of the canonical object.

```

1: for  $i = 1$  to  $K$  do
2:   for  $j = 1$  to  $K, j \neq i$  do
3:      $C_{i,j} = \frac{1}{|X^{(j)}|} \sum_{k=1}^{|X^{(j)}|} \min_{l=1}^{|X^{(i)}|} \|X_k^{(j)} - X_l^{(i)}\|_2^2$ 
4:   end for
5: end for
6: for  $i = 1$  to  $K$  do
7:    $C_i = \sum_{j=1, j \neq i}^K C_{i,j}$ 
8: end for
9: return  $\arg \min_{i=1}^K C_i$ 

```

421 A.3 Gram-Schmidt Orthogonalization

422 We compute a rotation matrix from two 3D vectors using Algorithm 5 [38].

¹<https://github.com/mikedh/trimesh>

Algorithm 5 Gram-Schmidt Orthogonalization

Input: 3D vectors u and v .

Output: Rotation matrix.

- 1: $u' = u / \|u\|$
 - 2: $v' = \frac{v - (u' \cdot v)u'}{\|v - (u' \cdot v)u'\|}$
 - 3: $w' = u' \times v'$
 - 4: **return** Stack(u', v', w')
-

423 A.4 Shape and Pose Inference Details

424 The point clouds $Y \in \mathbb{R}^{n \times 3}$ starts in its canonical form with the latent shape v equal to zero. We set
425 the initial scale s to one, translation t to zero and rotation \hat{R} to identity,

$$v = \underbrace{(0 \quad 0 \quad \dots \quad 0)}_d, \quad s = (1 \quad 1 \quad 1), \quad t = (0 \quad 0 \quad 0), \quad \hat{R} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \end{pmatrix}. \quad (9)$$

426 \hat{R} is then transformed into $R \in \text{SO}(3)$ using Algorithm 5. We minimize \mathcal{L} with respect to v, s, t
427 and \hat{R} using the Adam optimizer [36] with learning rate 10^{-2} for 100 steps. We set $\beta = 10^{-2}$. We
428 found the optimization process is prone to getting stuck in local minima; e.g., instead of aligning
429 the handle of the decoded mug with the observed point cloud, the optimizer might change the shape
430 of the decoded mug to hide its handle. Hence, we restart the process with many different random
431 initial rotations and pick the solution with the lowest loss function. Further, we randomly subsample
432 Y to 1k points at each gradient descent step – this allows us to run 12 random starting orientations
433 at once on an NVIDIA RTX 2080Ti GPU.

434 A.5 Using Multiple Demonstrations

435 Our method transfers grasps and placements from a single demonstration, but in our simulated ex-
436 periment, we have access to multiple demonstrations. We implement a simple heuristic for choosing
437 the demonstration that fits our method the best: we make a prediction of the relational object place-
438 ment from the initial state of each demonstration and select the demonstration where our prediction
439 is closest to the demonstrated placement. The intuition is that we are choosing the demonstration
440 where our method was able to warp the objects with the highest accuracy (leading to the best place-
441 ment prediction). This is especially useful in filtering out demonstrations with strangely shaped
442 objects.

443 B Experiment Details

444 B.1 Object re-arrangement on a physical robot

445 We use a UR5 robotic arm with a Robotiq gripper. We capture the point cloud using three RealSense
446 D455 camera with extrinsics calibrated to the robot. For motion planning, we use MoveIt with
447 ROS1. To segment the objects, we use DBSCAN to cluster the point clouds and simple heuristics
448 (e.g. height, width) to detect the object class.

449 B.2 Grasp prediction in the wild

450 We use a single RealSense D435 RGB-D camera. Our goal is to be able to demonstrate any task
451 in the real world without having to re-train our perception pipeline. Therefore, we chose an open-
452 vocabulary object detection model Detic [42], which is able to detect object based on natural lan-
453 guage descriptions. We used the following classes: "cup", "bowl", "mug", "bottle", "cardboard",
454 "box", "Tripod", "Baseball bat", "Lamp", "Mug Rack", "Plate", "Toaster" and "Spoon". We use

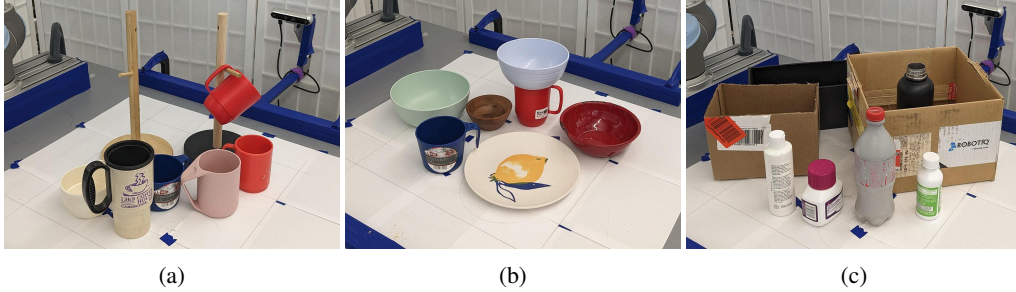


Figure 8: Objects used for the real-world tasks: (a) mug on tree, (b) bowl (or plate) on mug and (c) bottle in box. We use a single pair of objects to generate demonstrations and test on novel objects.

the predicted bounding boxes from Detic to condition a Segment Anything model [43] to get accurate class-agnostic segmentation masks. Both Detic² and Segment Anything³ come with several pre-trained models and we used the largest available. Finally, we select the pixels within each segmentation mask and use the depth information from our depth camera to create a per-object point cloud. We use DBSCAN to cluster the point cloud and filter out outlier points. Then, we perform mesh warping and interaction warping to predict object meshes and grasps.

Previously, we experimented with Mask R-CNN [44] and Mask2Former [45] trained on standard segmentation datasets, such as COCO [46] and ADE20k [47]. We found that these dataset lack the wide range of object classes we would see in a household environment and that the trained models struggle with out-of-distribution viewing angles, such as looking from a steep top-down angle. We also experimented with an open-vocabulary object detection model OWL-ViT [48] and found it to be sensitive to scene clutter and the viewing angle.

C Additional Results

Training and inference times: We measure the training and inference times of TAX-Pose, R-NDF and IW (Table 3). Both R-NDF and IW take tens of seconds to either perceive the environment or to predict an action. This is because both of these methods use gradient descent with many random restarts for inference. On the other hand, TAX-Pose performs inference in a fraction of second but requires around 16 hours of training for each task. Neither R-NDF nor IW require task-specific training. We do not include the time it takes to perform pre-training for each class of objects, which is required by all three methods, because we used checkpoints provided by the authors of TAX-Pose and R-NDF.

Additional real-world grasp predictions: We include additional examples of real-world object segmentation, mesh prediction and grasp prediction in Figure 9.

²<https://github.com/facebookresearch/Detic>

³<https://github.com/facebookresearch/segment-anything>

Method	Training	Perception	Grasp prediction	Placement prediction
TAX-Pose [2]	16.5 ± 1.3 h	-	0.02 ± 0.01 s	0.02 ± 0.01 s
R-NDF [17]	-	-	21.4 ± 0.5 s	42.5 ± 1.8 s
IW (Ours)	-	29.6 ± 0.2 s	0.01 ± 0.01 s	0.003 ± 0.004 s

Table 3: Approximate training and inference times for our method and baselines measured over five trials. R-NDF and IW do not have an explicit training phase, as they use demonstrations nonparametrically during inference. Only IW has a perception step that is separate from the action prediction step. We do not include the time it takes to capture a point cloud or to move the robot. Training and inference times were measured on a system with a single NVIDIA RTX 2080Ti GPU and an Intel i7-9700K CPU.

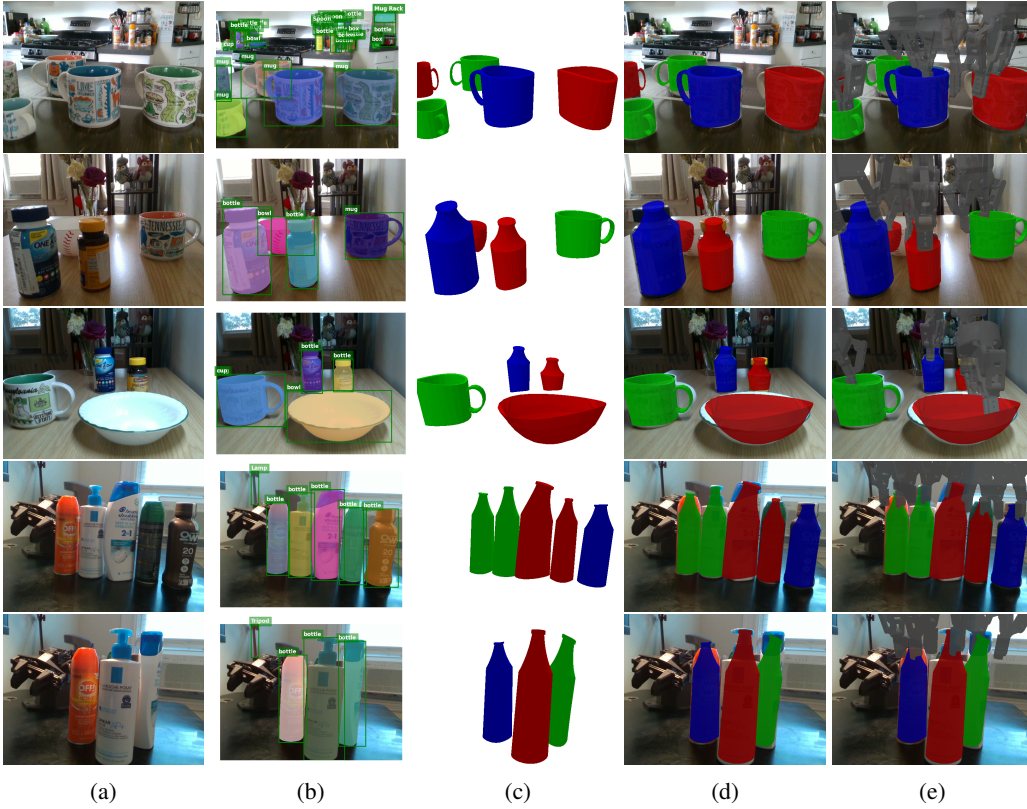


Figure 9: Additional examples, please see Figure 7.

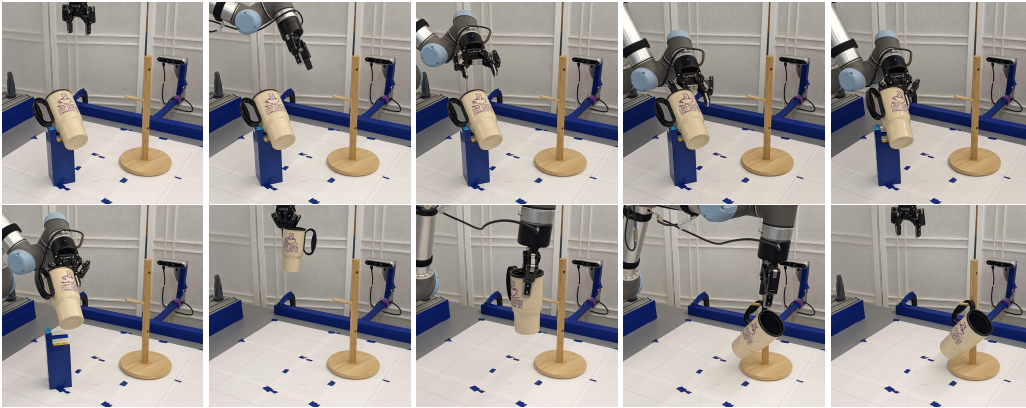


Figure 10: Example of mug on tree episode.

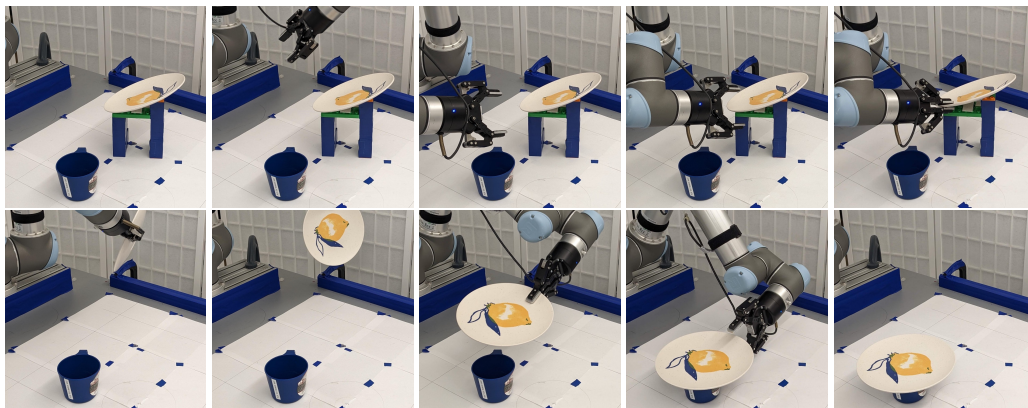


Figure 11: Example of bowl/plate on mug episode.

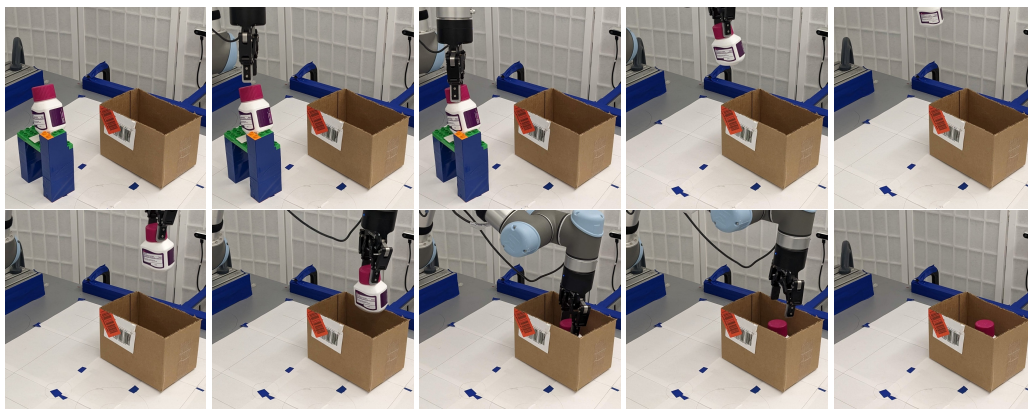


Figure 12: Example of bottle in box episode.