

Reactive In-Air Clothing Manipulation with Confidence-Aware Dense Correspondence and Visuotactile Affordance

Anonymous Author(s)

Affiliation

Address

email

Abstract: Manipulating clothing is challenging due to their complex, variable configurations and frequent self-occlusion. While prior systems often rely on flattening garments, humans routinely identify keypoints in highly crumpled and suspended states. We present a novel, task-agnostic, visuotactile framework that operates directly on crumpled clothing—including in-air configurations that have not been addressed before. Our approach combines global visual perception with local tactile feedback to enable robust, reactive manipulation. We train dense visual descriptors on a custom simulated dataset using a distributional loss that captures cloth symmetries and generates correspondence confidence estimates. These estimates guide a reactive state machine that dynamically selects between folding strategies based on perceptual uncertainty. In parallel, we train a visuotactile grasp affordance network using high-resolution tactile feedback to supervise grasp success. The same tactile classifier is used during execution for real-time grasp validation. Together, these components enable a reactive, task-agnostic framework for in-air garment manipulation, including folding and hanging tasks. Moreover, our dense descriptors serve as a versatile intermediate representation for other planning modalities, such as extracting grasp targets from human video demonstrations, paving the way for more generalizable and scalable garment manipulation.

Keywords: Deformable Object Manipulation, Dense Correspondence Learning, Confidence-Aware Planning, Visuotactile Perception

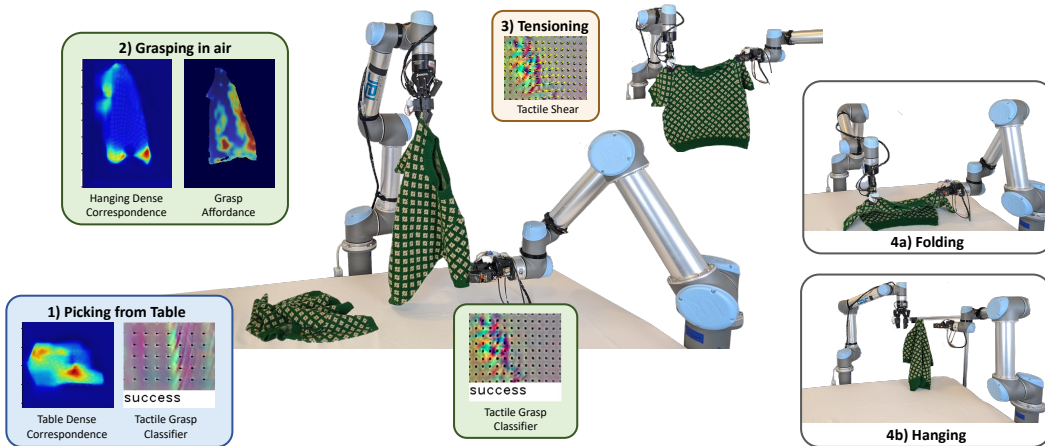


Figure 1: **Overview of visuotactile garment manipulation system.** Our framework integrates dense visual correspondence, visuotactile grasp affordance prediction, tactile grasp evaluation, and tactile tensioning for manipulating garments in crumpled configurations, both on a table-top and in-air. By leveraging a confidence-aware, reactive architecture and a task-agnostic representation, the system supports a variety of manipulation tasks—including folding and hanging.

21 1 Introduction

22 Deformable object manipulation remains a major challenge in robotics, since strategies developed
23 for rigid objects often fail to transfer. Deformable objects occupy infinite-dimensional configuration
24 spaces and exhibit high model uncertainty, making accurate state estimation and dynamics prediction
25 difficult. Although simulation-based models exist, they are typically computationally intensive and
26 insufficiently inaccurate for real-time control. In this work, we focus on garment manipulation,
27 where real-world complexities—such as self-occlusion, intra-class variation, and diverse material
28 dynamics—further complicate perception and control.

29 Existing approaches typically fall into two extremes: full-state estimation, which is expensive, or
30 task-specific grasp predictors, which lack generalizability. To bridge this gap, we propose a pose-
31 and instance-agnostic, confidence-aware representation using dense visual descriptors that estab-
32 lishes pixel-wise correspondences between crumpled garments and canonical flat configurations.
33 Trained on highly deformed states of detailed simulated shirts, our model can directly identify corre-
34 spondences for shirts crumpled on a table and suspended in the air—a setting that, to our knowledge,
35 has not been previously addressed.

36 Instead of the traditional contrastive loss, we use a distributional loss that models garment symme-
37 tries and produces confidence estimates for each correspondence. These confidence scores inform
38 whether a keypoint should be grasped or deferred, which is critical for operating under severe occlu-
39 sion. We integrate this representation into a visuotactile manipulation system, using high-resolution
40 tactile sensing to (1) supervise grasp affordance learning, (2) validate grasp success during execu-
41 tion, and (3) enable closed-loop tensioning during folding. These components work together within
42 a reactive framework that adapts folding and hanging strategies to garments of varying geometries,
43 without requiring full-state estimation or flattening.

44 We make the following key technical contributions:

- 45 • **Parametrizable Simulator:** A custom simulator with realistic hem features and parameterized
46 variations to enable correspondences across different shirt geometries.
- 47 • **Dense Representation:** Pixel-wise correspondences across challenging states using a distribu-
48 tional loss to capture symmetries and provide confidence estimates.
- 49 • **Visuotactile Affordance:** Grasp affordance network trained in simulation and fine-tuned using
50 tactile supervision.
- 51 • **Cloth Manipulation System:** A reactive visuotactile framework combining dense correspon-
52 dences, affordances, and tactile sensing for confidence-aware in-air folding and hanging.

53 2 Related Works

54 Most previous cloth manipulation work focuses on task-specific pipelines, including flattening [1, 2],
55 folding [3, 4, 5], dressing [6, 7], and recently hanging [8, 9, 10, 11]. These systems typically use
56 incremental pick-and-place motions against a table [12, 13, 5, 14], and many focus on rectangular
57 cloth, rather than garments.

58 Learning-based approaches can be quite successful at specific tasks. Labeling a real-world de-
59 formable object dataset is challenging [15, 9], so most learning works are trained in simulation.
60 However, the sim2real gap remains a challenge—we address this for our grasp affordance network
61 by extending [16], fine-tuning using tactile classifiers to determine grasp success on the robot. Be-
62 havior cloning approaches [8] have shown impressive results on tasks like tying shoelaces and hang-
63 ing shirts, but require thousands of expert teleoperated demonstrations per task. In contrast, our
64 system enables one- or few-shot generalization abilities and can reuse a shared object-centric repre-
65 sentation across tasks.

66 **Perception and Representation** Early cloth manipulation work relies on corner detection or ridge
67 detection [17] to determine grasp points [18]. However, finding other more specific local features

often requires first flattening the cloth [12, 19, 14, 1] or hanging it from specific grasp points [4, 20, 3] to avoid self-occlusion. Some works determine the global state of the cloth [21, 22, 23], but full-state inference is computationally expensive. In contrast, we use dense pixel-wise correspondences to directly localize task-relevant points in both crumpled table-top and in-air configurations.

Dense Descriptors Dense visual descriptors have been used to learn pixel-level correspondences across object views [24, 25]. Florence et al. [26] introduce dense object descriptors for task-agnostic manipulation, with follow-up work applying them to deformable objects [5, 27, 28]. Prior cloth-specific applications use contrastive loss [5, 28], but Ganapathi et al. [29] use multimodal distributional loss [30] to model symmetry and uncertainty on ropes and square cloths. We extend this to garments, training on highly crumpled configurations and enabling in-air correspondence prediction—a capability not previously addressed. Our approach further differs from garment manipulation in [28] because of our use of reactive control, made possible by confidence-aware descriptors and tactile feedback. We also demonstrate that our dense descriptors can act as an intermediate representation for different planning modalities. For example, Huang et al. [31] uses DinoV2 [32] and a vision-language model to determine constraints; our descriptors could find keypoint candidates to better support manipulation in more crumpled states.

3 Methods

3.1 Dataset Generation in Simulation

We use Blender 4.2 [33] to simulate a wide variety of shirt geometries and deformations, generating a large RGB-D dataset (1500 scenes) for training. In addition to parameterizing the overall geometries, we use [34] to incorporate hems, stitches, and sewing seams into our shirts to mimic realistic garments, enhancing visual realism and providing key features helpful for correspondence. Our method incorporates these finer details while preserving consistent vertex indexing across shirts, enabling descriptors to align with a canonical template regardless of geometry, without relying on sparse skeleton keypoints as in [28]. Figure 2 shows some of the parameters and shirt configurations we randomize to generate our dataset.

Scene generation mimics real-world camera setups, with three cameras arranged radially around the hanging shirt, with added pose noise and varied lighting conditions to enhance dataset diversity. For each hanging scene, a shirt is hung from a random mesh point and the world coordinates and pixel locations of the deformed mesh vertices are saved. For each table scene, a randomly positioned flat shirt is repeatedly grasped from random points and repositioned multiple times. This setup captures rich, diverse data across garment shapes, crumpled configurations (hanging and table), and visual contexts, enabling robust correspondence learning between different poses and shirt instances. See Appendix for further simulation details.

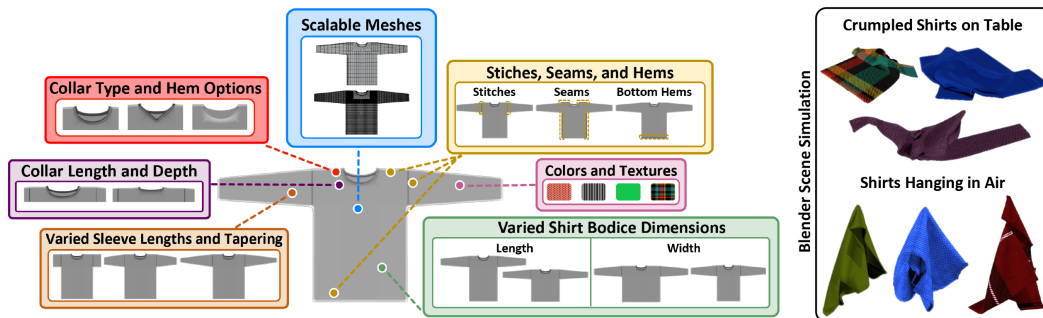


Figure 2: **Generating a simulated shirt dataset.** Blender 4.2 is used to simulate deformed shirts. Our animation pipeline allows flexibility in shirt geometries with the addition of realistic, key features like seams and hems often found on real shirts. A consistent vertex indexing across the shirt dataset is used, allowing alignment with a canonical template.

3.2 Dense Correspondence with Distributive Loss

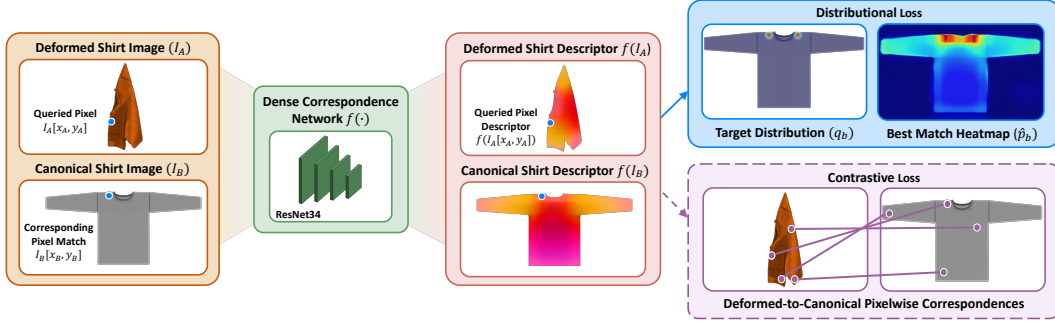


Figure 3: **Training dense correspondence in simulation.** Given two images I_a and I_b , and a matching relation $((x_a, y_a) \longleftrightarrow \{(x_b, y_b), (x'_b, y'_b)\})$, we train a CNN model f to compute dense object descriptors. When supervising with distributional loss, we define a multimodal Gaussian target distribution q_b with symmetrical modes over pixels corresponding to the queried point. We compute the probability distribution estimation \hat{p}_{b_i} over image I_b using $f(I_a)[x_a, y_a]$ and $f(I_b)$. Training minimizes the KL divergence between q_b and \hat{p}_{b_i} . In the contrastive loss case, the model learns to push discrete pixel matches closer together in pixel space and non-matches further apart.

We aim to learn dense pixel-wise correspondences between images of deformable objects in crumpled and flattened configurations. Given an RGB image $I \in \mathbb{R}^{W \times H \times 3}$, we define a mapping $f : \mathbb{R}^{W \times H \times 3} \rightarrow \mathbb{R}^{W \times H \times d}$ that assigns a d -dimensional descriptor to each pixel in I . This descriptor space allows correspondences to be established by comparing descriptors across images.

Contrastive Loss Contrastive methods, as used by [26, 5, 28], supervise this mapping by sampling pairs of matching and non-matching pixels across images. For a query pixel $u_a = (x_a, y_a)$ in image I_a and a candidate pixel $u_b = (x_b, y_b)$ in image I_b , the descriptor distance $D(I_a, u_a, I_b, u_b) = \|f(I_a)(u_a) - f(I_b)(u_b)\|_2$ is minimized for matching pairs and pushed apart for non-matching pairs. This enforces one-to-one correspondences but struggles with ambiguities caused by symmetries or occlusions, which are common in deformable objects. Symmetric Pixel-wise Contrastive Loss (SPCL) [29] extends this approach to support symmetric correspondences, allowing multiple valid matches per query pixel. However, they found the results to be unstable, and the discrete matches resulted in discontinuity issues. We will compare our network to these contrastive baselines.

Distributional Loss To address these limitations, we adopt the distributional formulation from [29], which directly models uncertainty over correspondences. Instead of supervising individual descriptor pairs, the network predicts a full probability distribution over possible matches. Specifically, we define an estimator $\hat{p}_b(x_i, y_j | I_a, I_b, x_a, y_a)$ that outputs the probability that each pixel $(x_i, y_j) \in I_b$ corresponds to a given query pixel $(x_a, y_a) \in I_a$. This estimator is defined as:

$$\hat{p}_b(x_i, y_j | I_a, I_b, x_a, y_a) = \frac{\exp \|f(I_a)[x_a, y_a] - f(I_b)[x_i, y_j]\|_2^2}{\sum_{i', j'} \exp \|f(I_a)[x_a, y_a] - f(I_b)[x_{i'}, y_{j'}]\|_2^2} \quad \forall (x_i, y_j) \in I_b \quad (1)$$

The target distribution q_b is a multimodal isotropic Gaussian defined over I_b , with standard deviation σ and modes centered at the ground-truth correspondence pixels, allowing the network to represent multiple valid matches and capture ambiguities from symmetry.

The descriptor mapping f is implemented using ResNet34. The network is optimized by minimizing the Kullback-Leibler (KL) divergence between the predicted distribution \hat{p}_{b_i} and the target distribution q_{b_i} for each query pixel. Here, \hat{p}_{b_i} is the predicted correspondence distribution over I_b for the i -th query pixel (computed using Equation 1), and q_{b_i} is the corresponding target distribution. Figure 3 shows a training example. At each iteration, we choose an image of a randomized crumpled shirt and compare it to the canonical one. We query 50 randomly sampled points on the crumpled shirt per iteration.

Note that I_b is always the canonical shirt image, meaning that we compute both the target and estimated distributions over the canonical shirt. A smooth Gaussian target distribution works over the canonical shirt because it does not have occlusions and distortions of the crumpled shirt. Defining the target distribution over the crumpled shirt would be useful for training the network in both directions, but is unfeasible in this framework.

3.3 Visuotactile Grasp Affordance

Training a general garment grasp affordance network is more challenging than for simpler deformable objects like towels. In [16], the network was fine-tuned on a single towel with consistent material properties and dynamics. However, in this case, affordance must generalize across a wide range of geometries and material rigidities. As in [16], we only use side grasps to reduce computational complexity. While grasp classifiers are trained for both grippers (as required by the larger system), affordance training is performed only for right-arm grasps, with left-arm affordance approximated by horizontally flipping inputs and outputs.

Tactile Classifier To assess grasp quality, we train tactile classifiers to distinguish between successful grasps, grasps with too little fabric (which are prone to slip), and grasps with excess layers (indicating more fabric than intended). We concatenate five evenly-spaced tactile depth images from the grasp attempt as input to our network. Our tactile datasets includes 350 grasps across approximately 20 shirts, with limited augmentations (two per input).

Training Affordance in Simulation We use the same U-Net [35] architecture as [16] for affordance prediction. The input to the network is a depth image of the hanging garment, and the output is an affordance heatmap over the image. Ground-truth affordance labels are computed per pixel via geometric analysis, leveraging full access to the cloth state in simulation. Specifically, each pixel is labeled based on gripper reachability, collision avoidance, and the number of fabric layers inside the gripper (restricted to two or fewer). These criteria are all explicitly checked in simulation, but the tactile classifier implicitly verifies these qualities on the robot. The simulated dataset consists of 300 unique cloth configurations, each rotated in increments of 30° , yielding a total of 3,600 images.

Fine-tuning on the Robot We collect 8,500 grasp points for real-world fine-tuning to capture the greater variety of shirt dynamics and configurations compared to the simulated environment. Fine-tuning can easily overfit the real grasp dataset because the loss only applies to one pixel at a time. Furthermore, the tactile classifier cannot reliably determine whether the grasped region corresponds to the intended visual target. As a result, non-reachable pixels can yield positive tactile signals due to inadvertently grasping cloth in front of the target. To help address these challenges, our loss includes neighboring pixels to broaden supervision, along with regularization terms such as spatial smoothness penalties, simulation consistency constraints, and weight decay.

3.4 System Setup

Our bimanual system consists of two UR5 robots, both equipped with parallel-jaw grippers mounted with GelSight Wedge tactile sensors [36]. A Kinect Azure camera is used to capture RGB-D images.

3.5 In-Air Garment Manipulation

Folding with Confidence-based State Machine Unlike prior garment folding approaches that rely on fixed canonical keypoints [5, 28] for folding on a table, our system enables reactive in-air folding by dynamically selecting grasp points based on real-time confidence estimates and recovering from failures using tactile reactivity. The system starts by picking the shirt up from the table (looking for high-confidence correspondence regions), and all subsequent grasps are performed in air.

At each grasp attempt, the robot can query from three canonical regions (shoulder, sleeve, bottom) using our distributional dense correspondence network to generate confidence-weighted heatmaps. A grasp is executed only if both the correspondence confidence and grasp affordance (for hanging grasps) exceed predefined thresholds. Otherwise, the robot rotates the garment by 30° and re-

178 evaluates, ensuring robust grasp point selection across four folding strategies (shoulder-to-shoulder,
179 bottom-to-bottom, sleeve-to-sleeve, sleeve-to-bottom) (See Appendix for details).

180 Grasp success is validated by tactile sensing (confirming fabric contact). If a grasp fails, the robot
181 rotates and retries without releasing the garment. We use vision to ensure that the cloth is still in
182 grip after moving the grippers. If no pixel meets the threshold requirements, the robot grasps the
183 lowest available high affordance point to change configurations and encourage the cloth to unfurl.
184 Once two confident grasp points are secured, the robot tensions the shirt (detecting shear via marker
185 tracking on tactile sensors) and performs the rest of the fold motions open-loop.

186 **Hanging** We demonstrate hanging by picking collar or shoulder from the table and in the air. After
187 securing both grasps, the robot moves open-loop to a peg. Hanging success is evaluated by grasp
188 regions and whether the cloth stays on the peg.

189 4 Results

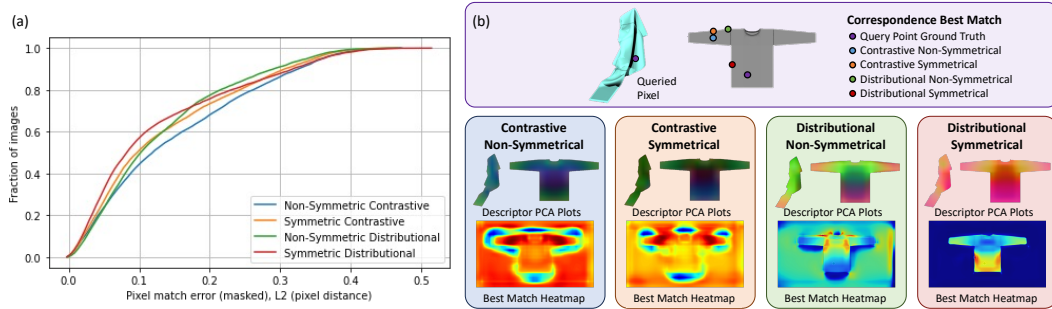


Figure 4: (a) **Cumulative pixel match error curves comparing contrastive and distributional training, with and without symmetric supervision along with (b) illustrative example.** The networks were trained on a combined dataset of hanging and table shirts and (a) shows performance on an unseen hanging test set. Higher curves indicate better performance. For each network, we show the predicted best pixel match for a queried point on a crumpled simulated shirt (b). We also provide PCA visualizations of the dense descriptors in both the canonical and crumpled states, alongside the corresponding match heatmaps. Note that contrastive heatmaps are normalized between 0 and 1 for visualization, while distributional heatmaps represent true correspondence probabilities.

190 **Dense Correspondence** Most dense descriptor methods use contrastive one-to-one training [26, 5,
191 28], which fails to capture symmetries or spatial relationships beyond binary matches. Quantitative
192 results (Fig. 5) show similar cumulative pixel errors between contrastive and distributional models,
193 but distributional models consistently outperform contrastive ones across nearly all error thresholds.
194 Qualitatively, contrastive loss struggles ambiguous structures, often collapsing descriptors along the
195 entire sleeve or confusing sleeves with the shirt bottom (as seen in PCA visualizations). In contrast,
196 distributional loss supervises the model to predict a full probability distribution, enforcing spatial
197 consistency. Explicit symmetry supervision further improves performance (Fig. 5), especially at
198 low error thresholds, by encouraging multimodal correspondences in symmetric regions.

199 We found that including occlusions during training did not significantly affect performance in simu-
200 lation, but helped improve performance on real data, likely due to masking artifacts. More detailed
201 analysis of network parameters can be found in the Appendix.

202 On real robot hanging images, we evaluate our network by defining classification zones on the
203 canonical shirt (see Appendix). When querying points from a crumpled hanging shirt (forward
204 direction), the best hanging-only network classified the correct region 73.3% of the time, while the
205 best combined network (trained on both table and hanging data) achieved 62.2% accuracy, while
206 exhibiting lower overall confidence. Applying a confidence threshold, the combined network made
207 correct, confidence-aware decisions (avoiding incorrect labels) 68.9% of the time. In the inverse
208 direction (querying from the canonical shirt), the combined network correctly identified the region
209 41.7% of the time and made safe, confidence-aware decisions 70.8% of the time. Some canonical

points were occluded in the crumpled image, making low confidence the correct outcome for these cases. On table scenes, the correct correspondence region was identified 70% of the time, and a safe decision—either correct or low-confidence—was made 80% of the time in 20 trials.

Visuotactile Grasp Affordance Our tactile grasp classifier achieves 99.7% accuracy on the right arm (used for tactile supervision) and 98.8% on the left. Thin, flat shirts are the most challenging to classify. To evaluate affordance prediction, we collect 125 human-labeled grasp points where each point appeared potentially graspable to a human observer. We compare our fine-tuned affordance network against two baselines: (1) Sim2Real, trained in simulation and directly deployed, and (2) Real2Real, trained solely on robot data. Networks are evaluated offline using precision@k [37], a metric suitable for our unbalanced test set that avoids the need for a fixed threshold. We report precision@80, corresponding to the 80 successful grasps among the 125 test points. The results are 71.3% for Sim2Real, 75.0% for Real2Real, and 76.3% for our fine-tuned network. Sim2Real performs worst due to discrepancies between simulated and real-world dynamics. While the fine-tuned and Real2Real networks achieve similar precision, qualitative analysis shows that Real2Real tends to be overconfident in incorrect predictions, particularly in less ambiguous cases not well-represented in the test set (see Appendix).

Combined System We evaluate grasping performance across four garment regions—sleeve, bottom, shoulder, and collar—using two networks: one trained solely on hanging data and another on a combined table and hanging dataset. For each category, we perform 10 grasp attempts per network, recording outcomes as success, failure, or below confidence threshold. Failures are further categorized as correspondence errors or affordance errors. In this experiment, we place the shirt in configurations where we expect graspable regions to emerge after rotation. Table 1 summarizes rates for overall success, correspondence success (excluding bad affordance grasps), low-confidence rates, and total failure rates for each network and region.

The collar region consistently achieves higher confidence and success rates, likely due to its distinctive geometry. In contrast, the bottom region has the lowest confidence rates, reflecting its visual ambiguity and the increased difficulty of finding good affordance grasps from folding in on itself. The hanging network performs marginally better overall, but the combined network adds critical flexibility by supporting table grasps. Importantly, during folding, we query three candidate grasp points for the initial grasp, requiring confidence in only one to proceed. Subsequent grasps occur in easier, more unfurled configurations.

Category	Successful Grasp (%)		Corr. Success (%)		Low Conf. (%)		Failed Grasp (%)	
	Hang	Comb	Hang	Comb	Hang	Comb	Hang	Comb
Sleeve	60	40	80	60	10	10	30	50
Bottom	40	10	90	90	40	80	20	10
Shoulder	40	60	100	100	60	20	0	20
Collar	80	80	90	90	0	0	20	20

Table 1: Grasping results using dense correspondence and grasp affordance across shirt categories for hanging and combined (hanging + table) dataset networks. Low-confidence outcomes, where the shirt completes a full rotation without finding a grasp point, are not counted as successful or failed grasps. They are still included when calculating correspondence success, since both networks are trained to be confidence-aware. Failed grasps are categorized as either correspondence or affordance failures. Correspondence success rates exclude grasps that failed due to bad affordance predictions.

We found that our confidence-aware state machine was able to grasp viable folding points in 6 out of 10 trials. Irrecoverable failure modes included correspondence failures, grabbing too much fabric, and grabbing diagonally across the shirt for sleeve-end grasps (despite masking out lowest points, see Appendix). Cloth slipping out was an occasional issue, but the system is able to recover. Our hanging system was successful in 7 out of 10 trials with all failures due to correspondence.

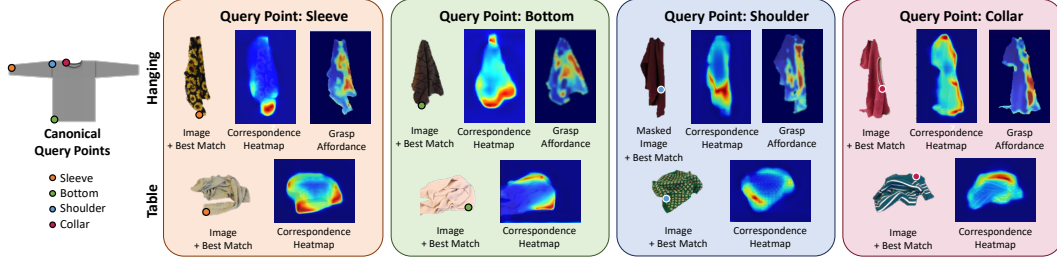


Figure 5: **Correspondence and affordance heatmaps for real images.** We show examples for both hanging and table configurations, with correspondence probability maps for four query types: sleeve, shoulder, collar, and bottom. For hanging images, we also show the grasp affordance heatmap. In the robot system, grasp points are selected where both correspondences and affordance exceed predefined confidence thresholds. Note that while training queries points on the crumpled shirt, the robot queries points on the canonical image.

5 Conclusion

We present a reactive visuotactile system for garment manipulation that integrates dense visual correspondence, visuotactile grasp affordance, confidence-aware planning, and tactile feedback. Unlike prior work constrained to table-top picking or reliant on flattening, our system supports in-air garment manipulation directly from crumpled states, guided by dense correspondences—a capability not previously demonstrated in the field. This enables more flexible, human-like manipulation.

A core insight of our work is the importance of confidence-driven reactivity: by deferring low-confidence actions and using tactile sensing for validation and correction, the system maintains robustness under severe occlusion and uncertainty. This closed-loop approach bridges the gap between global visual context and local contact feedback, enabling reliable control even when full object geometry is not observable.

Beyond task execution, our dense, confidence-aware representation serves as a generalizable intermediate layer for higher-level planning frameworks. It provides a foundation for extracting grasp targets from human video demonstrations (Fig 6, See Appendix for details), and has the potential to interface with vision-language models [31] or symbolic planners. These directions open the door to scalable, semantically-informed manipulation systems capable of adapting across garments, tasks, and contexts.

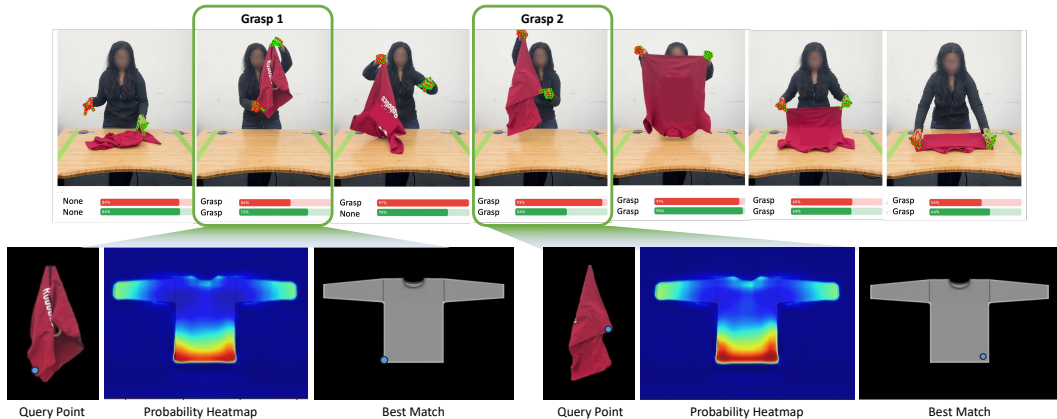


Figure 6: **Extracting grasp points from human video demonstrations.** We track hand gestures throughout the video to identify key moments. For each key frame, we use the tracked hand position to define a query point and retrieve the corresponding location on the canonical shirt using our dense correspondence model. This approach enables folding demonstrations to be interpreted as robot-executable instructions via our dense visual representation.

6 Limitations

While our system demonstrates strong potential for in-air garment manipulation, several areas present opportunities for further development. First, the generalizability of the dense correspondence network is limited by the features available in simulation. Although we incorporated realistic details such as seams, hems, and varied necklines, other common garment features—like hoods, buttons, zippers, and mixed patterns—are not yet modeled. Some of these could be added in future dataset expansions, while others may require advances in simulation tools. On out-of-distribution shirts (see Appendix), the network still captures general structure, but with lower confidence.

Second, we are able to achieve this performance with a single camera and exclusively side approach grasps, but expanding to additional viewpoints and enabling more grasp approach angles could improve coverage to access more high correspondence regions. Incorporating temporal information could further enable the system to track keypoints as they become accessible, supporting more flexible planning.

Finally, although the system is confidence-aware, the network occasionally overestimates its certainty in challenging configurations. We experimented with auxiliary confidence prediction and KL-divergence metrics, but these did not significantly improve failure detection. Improving uncertainty estimation remains an important direction for future work.

References

- [1] H. Ha and S. Song. Flingbot: The unreasonable effectiveness of dynamic manipulation for cloth unfolding. In *Conference on Robot Learning*, pages 24–33. PMLR, 2022.
- [2] A. Canberk, C. Chi, H. Ha, B. Burchfiel, E. Cousineau, S. Feng, and S. Song. Cloth funnels: Canonicalized-alignment for multi-purpose garment manipulation, 2022. URL <https://arxiv.org/abs/2210.09347>.
- [3] J. Maitin-Shepard, M. Cusumano-Towner, J. Lei, and P. Abbeel. Cloth grasp point detection based on multiple-view geometric cues with application to robotic towel folding. In *2010 IEEE International Conference on Robotics and Automation*, pages 2308–2315. IEEE, 2010.
- [4] A. Doumanoglou, A. Kargakos, T.-K. Kim, and S. Malassiotis. Autonomous active recognition and unfolding of clothes using random decision forests and probabilistic planning. In *2014 IEEE international conference on robotics and automation (ICRA)*, pages 987–993. IEEE, 2014.
- [5] A. Ganapathi, P. Sundaresan, B. Thananjeyan, A. Balakrishna, D. Seita, J. Grannen, M. Hwang, R. Hoque, J. E. Gonzalez, N. Jamali, K. Yamane, S. Iba, and K. Goldberg. Learning to smooth and fold real fabric using dense object descriptors trained on synthetic color images. *arXiv*, 2020.
- [6] F. Zhang and Y. Demiris. Learning garment manipulation policies toward robot-assisted dressing. *Science robotics*, 7(65):eabm6010, 2022.
- [7] Z. Sun, Y. Wang, D. Held, and Z. Erickson. Force-constrained visual policy: Safe robot-assisted dressing via multi-modal sensing. *IEEE Robotics and Automation Letters*, 2024.
- [8] T. Z. Zhao, J. Tompson, D. Driess, P. Florence, K. Ghasemipour, C. Finn, and A. Wahid. Aloha unleashed: A simple recipe for robot dexterity. *arXiv preprint arXiv:2410.13126*, 2024.
- [9] W. Chen, D. Lee, D. Chappell, and N. Rojas. Learning to grasp clothing structural regions for garment manipulation tasks. In *2023 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 4889–4895. IEEE, 2023.
- [10] Y. Chen, S. Wei, B. Xiao, J. Lyu, J. Chen, F. Zhu, and H. Wang. Robohanger: Learning generalizable robotic hanger insertion for diverse garments. *arXiv preprint arXiv:2412.01083*, 2024.
- [11] W. Chen, K. Li, D. Lee, X. Chen, R. Zong, and P. Kormushev. Graphgarment: Learning garment dynamics for bimanual cloth manipulation tasks. *arXiv preprint arXiv:2503.05817*, 2025.
- [12] Y. Wu, W. Yan, T. Kurutach, L. Pinto, and P. Abbeel. Learning to manipulate deformable objects without demonstrations, 2019.
- [13] R. Hoque, D. Seita, A. Balakrishna, A. Ganapathi, A. K. Tanwani, N. Jamali, K. Yamane, S. Iba, and K. Goldberg. VisuoSpatial foresight for multi-step, multi-task fabric manipulation, 2020.
- [14] X. Lin, Y. Wang, Z. Huang, and D. Held. Learning visible connectivity dynamics for cloth smoothing. In *Conference on Robot Learning*, pages 256–266. PMLR, 2022.
- [15] J. Qian, T. Weng, L. Zhang, B. Okorn, and D. Held. Cloth region segmentation for robust grasp selection. In *IEEE International Conference on Intelligent Robots and Systems*, 2020. ISBN 9781728162126. doi:10.1109/IROS45743.2020.9341121.
- [16] N. Sunil, S. Wang, Y. She, E. Adelson, and A. R. Garcia. Visuotactile affordances for cloth manipulation with local control. In *6th Annual Conference on Robot Learning*, 2022. URL <https://openreview.net/forum?id=s6NEzqZKaP->.

- [17] K. Yamazaki, K. Nagahama, and M. Inaba. Daily clothes observation from visible surfaces based on wrinkle and cloth-overlap detection. In *MVA*, pages 275–278, 2011.
- [18] B. Willimon, S. Birchfield, and I. Walker. Model for unfolding laundry using interactive perception. In *IEEE International Conference on Intelligent Robots and Systems*, 2011. ISBN 9781612844541. doi:10.1109/IROS.2011.6048796.
- [19] R. Hoque, K. Shivakumar, S. Aeron, G. Deza, A. Ganapathi, A. Wong, J. Lee, A. Zeng, V. Vanhoucke, and K. Goldberg. Learning to fold real garments with one arm: A case study in cloud-based robotics research. *arXiv preprint arXiv:2204.10297*, 2022.
- [20] M. Cusumano-Towner, A. Singh, S. Miller, J. F. O’Brien, and P. Abbeel. Bringing clothing into desired configurations with limited perception. In *2011 IEEE international conference on robotics and automation*, pages 3893–3900. IEEE, 2011.
- [21] T. Tang, Y. Fan, H.-C. Lin, and M. Tomizuka. State estimation for deformable objects by point registration and dynamic simulation. In *Intelligent Robots and Systems (IROS), 2017 IEEE International Conference on*. IEEE, 2017.
- [22] C. Chi and D. Berenson. Occlusion-robust deformable object tracking without physics simulation. In *Intelligent Robots and Systems (IROS), 2019 IEEE International Conference on*. IEEE, 2019.
- [23] C. Chi and S. Song. Garmentnets: Category-level pose estimation for garments via canonical space shape completion. *CoRR*, abs/2104.05177, 2021. URL <https://arxiv.org/abs/2104.05177>.
- [24] C. B. Choy, J. Gwak, S. Savarese, and M. Chandraker. Universal correspondence network. *CoRR*, abs/1606.03558, 2016. URL <http://arxiv.org/abs/1606.03558>.
- [25] T. Schmidt, R. A. Newcombe, and D. Fox. Self-supervised visual descriptor learning for dense correspondence. *IEEE Robotics and Automation Letters*, 2017.
- [26] P. R. Florence, L. Manuelli, and R. Tedrake. Dense object nets: Learning dense visual object descriptors by and for robotic manipulation. *CoRR*, abs/1806.08756, 2018. URL <http://arxiv.org/abs/1806.08756>.
- [27] P. Sundaresan, J. Grannen, B. Thananjeyan, A. Balakrishna, M. Laskey, K. Stone, J. E. Gonzalez, and K. Goldberg. Learning rope manipulation policies using dense object descriptors trained on synthetic depth data. *CoRR*, abs/2003.01835, 2020. URL <https://arxiv.org/abs/2003.01835>.
- [28] R. Wu, H. Lu, Y. Wang, Y. Wang, and H. Dong. Unigarmentmanip: A unified framework for category-level garment manipulation via dense visual correspondence. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2024.
- [29] A. Ganapathi, P. Sundaresan, B. Thananjeyan, A. Balakrishna, D. Seita, R. Hoque, J. E. Gonzalez, and K. Goldberg. MMGSD: Multi-Modal Gaussian Shape Descriptors for Correspondence Matching in 1D and 2D Deformable Objects. In *Intelligent Robots and Systems (IROS), 2020 IEEE International Conference on*, 2020.
- [30] P. Florence. *Dense visual learning for robot manipulation*. PhD thesis, Massachusetts Institute of Technology, 01 2020.
- [31] W. Huang, C. Wang, Y. Li, R. Zhang, and L. Fei-Fei. Rekep: Spatio-temporal reasoning of relational keypoint constraints for robotic manipulation. *arXiv preprint arXiv:2409.01652*, 2024.

- 368 [32] M. Oquab, T. Darcet, T. Moutakanni, H. Vo, M. Szafraniec, V. Khalidov, P. Fernandez, D. Haz-
 369 iza, F. Massa, A. El-Nouby, et al. Dinov2: Learning robust visual features without supervision.
 370 *arXiv preprint arXiv:2304.07193*, 2023.
- 371 [33] Blender Online Community. Blender. <https://www.blender.org/>, 2025. Version 4.2.
- 372 [34] A. Albisser. Procedural cloth sewing toolbox for blender 4.2+. [https://](https://alexandrealbisser.gumroad.com/l/ProceduralSewingToolbox)
 373 alexandrealbisser.gumroad.com/l/ProceduralSewingToolbox, 2024. Software tool.
- 374 [35] O. Ronneberger, P. Fischer, and T. Brox. U-net: Convolutional networks for biomedical im-
 375 age segmentation. In *International Conference on Medical image computing and computer-*
 376 *assisted intervention*, pages 234–241. Springer, 2015.
- 377 [36] S. Wang, Y. She, B. Romero, and E. Adelson. Gelsight wedge: Measuring high-resolution
 378 3d contact geometry with a compact robot finger. In *2021 IEEE International Conference on*
 379 *Robotics and Automation (ICRA)*, pages 6468–6475. IEEE, 2021.
- 380 [37] M. Sanderson. *Test collection based evaluation of information retrieval systems*. Now Pub-
 381 lishers Inc, 2010.