

Foldsformer: Learning Sequential Multi-Step Cloth Manipulation with Space-Time Attention

Kai Mo¹, Chongkun Xia¹, Xueqian Wang¹, Yuhong Deng¹, Xuehai Gao² and Bin Liang³

Abstract—Sequential multi-step cloth manipulation is a challenging problem in robotic manipulation, requiring a robot to perceive the cloth state and plan a sequence of chained actions leading to the desired state. Most previous works address this problem in a goal-conditioned way, and goal observation must be given for each specific task and cloth configuration, which is not practical and efficient. Thus, we present a novel multi-step cloth manipulation planning framework named Foldsformer. Foldsformer can complete similar tasks with only a general demonstration and utilize a space-time attention mechanism to capture the instruction information behind this demonstration. We experimentally evaluate Foldsformer on four representative sequential multi-step manipulation tasks and show that Foldsformer significantly outperforms state-of-the-art approaches in simulation. Foldsformer can complete multi-step cloth manipulation tasks even when configurations of the cloth (e.g., size and pose) vary from configurations in the general demonstrations. Furthermore, our approach can be transferred from simulation to the real world without additional training or domain randomization. Despite training on rectangular clothes, we also show that our approach can generalize to unseen cloth shapes (T-shirts and shorts). Videos and source code are available at: <https://sites.google.com/view/foldsformer>.

Index Terms—Deep Learning in Grasping and Manipulation, Perception-Action Coupling

I. INTRODUCTION

CLOTH manipulation has wide applications in domestic, medical and industrial settings, such as laundry-folding [1], [2], surgery [3] and manufacturing [4]. Most of these tasks can be seen as sequential multi-step cloth manipulation, where sequential chained actions should be planned to achieve the desired cloth state. There are some challenges for sequential multi-step cloth manipulation: Unlike

Manuscript received: September, 2, 2022; Revised November, 18, 2022; Accepted December, 9, 2022. This paper was recommended for publication by Editor Markus Vincze upon evaluation of the Associate Editor and Reviewers' comments. This work was supported by the National Key R&D Program of China (2022YFB4701400/4701402), National Natural Science Foundation of China (No. U1813216, U21B6002, 62203260), the Shenzhen Science Fund for Distinguished Young Scholars (RCJC20210706091946001), Guangdong Young Talent with Scientific and Technological Innovation (2019TQ05Z111), China Postdoctoral Science Foundation (2022M711823). (Corresponding author: Chongkun Xia & Xueqian Wang.)

¹Kai Mo, Chongkun Xia, Xueqian Wang, and Yuhong Deng are with the Center for Intelligent Control and Telescience, Tsinghua Shenzhen International Graduate School, 518055 Shenzhen, China. (email: mok21@mails.tsinghua.edu.cn, xiachongkun@sz.tsinghua.edu.cn, wang.xq@sz.tsinghua.edu.cn, dengyh20@mails.tsinghua.edu.cn)

²Xuehai Gao is with the Research Institute of Tsinghua University in Shenzhen, 518055 Shenzhen, China. (email: gaohx@tsinghua-sz.org)

³Bin Liang is with the Navigation and Control Research Center, Department of Automation, Tsinghua University, 100084 Beijing, China. (email: bliang@tsinghua.edu.cn)

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

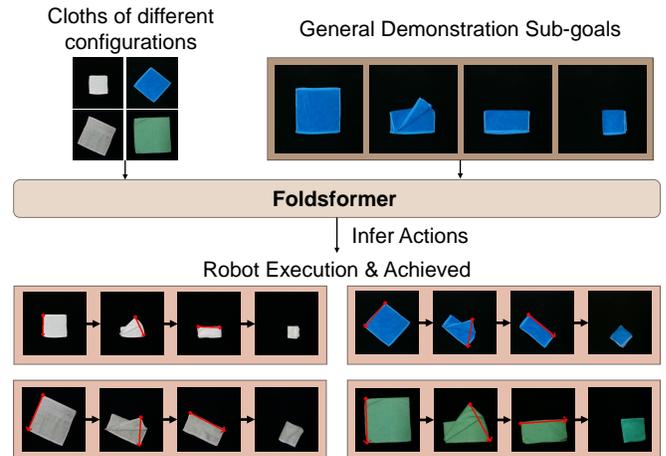


Fig. 1. **Overview.** Foldsformer can complete similar tasks even when the cloth configurations vary from that in a general demonstration.

rigid objects, cloth has infinite degrees of freedom, making it hard for state representation; The cloth dynamics is also complex [5], and slightly different interactions may lead to significantly different cloth behaviors; Further, a single counter-productive action may crumple the cloth and make it difficult to recover [6], which means completing a sequential multi-step cloth manipulation task requires only a particular set of chained actions.

For multi-step cloth manipulation, early works designed hard-coded policies with task-customized grippers and cloth with markers, making it hard to generalize these policies to other settings [1], [2]. Recently, learning-based approaches have shown great potential for cloth manipulation by tackling the difficulties in state estimation or modeling with neural networks. Some of these approaches learn task-specific sequential multi-step cloth manipulation policies from expert demonstrations. Although these approaches are demonstrated to be effective, different models should be trained to perform different tasks [7], [8], [9]. Others learn a goal-conditioned policy from random data and can achieve arbitrary cloth goal states. Lee *et al.* [10] specifies the goal state with a single image at test time. However, for a sequential multi-step cloth manipulation task, the cloth of the goal state can be highly self-occluded, and a single image fails to contain enough information about the goal. Weng *et al.* [11] specifies the goal state with a sequence of sub-goal images from an expert demonstration at test time. For the goal-conditioned policies mentioned above, the goal images must be given for each specific task and cloth configuration by human demonstrators, which is not practical and efficient. When configurations of the cloth vary from that in the goal, these policies' performances

greatly downgrade.

In this work, we address the problem mentioned above in sequential multi-step cloth manipulation. We propose a novel framework named Foldsformer that can complete similar tasks of different cloth configurations with a general demonstration. We design a novel encoder in Foldsformer based on the space-time attention mechanism. Given a general demonstration containing a sequence of sub-goals that describe a multi-step task, this encoder can capture the instruction information behind this demonstration that helps Foldsformer to perform similar tasks even when the cloth configuration varies (see Fig. 1). We first leverage large-scale pre-training from random actions and then transfer it to different tasks through only 100 expert demonstrations per task fine-tuning. Experimental results show that our approach greatly outperforms three previous state-of-the-art baselines [10], [11], [12]. We then demonstrate that our approach trained in simulation can be transferred to the real world without additional training or domain randomization. We also show that our learned model can successfully generalize to T-shirts and shorts despite being trained on rectangular clothes. In summary:

- We propose a novel framework named Foldsformer for sequential multi-step manipulation. Foldsformer can complete similar tasks on clothes of different configurations with a general demonstration.
- We propose a novel encoder architecture in Foldsformer based on space-time attention. This encoder enables a robot to learn instruction information behind a general demonstration.
- We conduct experiments to demonstrate that Foldsformer can be zero-shot transferred to the real world and generalize to unseen T-shirts and shorts despite being trained only on rectangular cloths.

II. RELATED WORK

A. Sequential Multi-Step Cloth Manipulation

There are two mainstream approaches for sequential multi-step cloth manipulation: analytical approaches and learning-based approaches. Analytical approaches are mainly task-specific, these approaches manually divide a task into several steps and then design scripted actions for each step. Doumanoglou *et al.* [1] and Bersch *et al.* [2] design complete pipelines for cloth folding. Yamakawa *et al.* [13] achieve dynamic folding of cloth with custom high-speed robot hands and high-speed sliders. Miller *et al.* [14] model cloth as polygons and complete folding tasks with a model-based optimization approach. However, these analytical approaches mainly use task-customized robot systems to adapt to their scripted actions. It is hard to generalize them to other settings.

Learning-based approaches can be divided into two categories: model-based and model-free. Model-based approaches learn a forward cloth dynamic model and use it for planning [15], [12], [16], [17]. However, these model-based approaches require time-consuming planning like Cross Entropy Method or Model Predictive Control to achieve the goal cloth state. As for model-free approaches, some approaches learn

task-specific cloth manipulation policies from expert demonstrations [7], [8], [9]. However, these task-specific policies fail to reuse information for different tasks efficiently. Some approaches learn a goal-conditioned policy that can achieve arbitrary unseen cloth states purely from random data. Lee *et al.* [10] learn a goal-conditioned pick-and-place policy with standard Deep Q-network and use predefined sets of fold angles and fold distances to ensure learning efficiency, but this discrete action space limits its manipulation ability. Weng *et al.* [11] propose a flow-based policy and divide a multi-step goal into a sequence of sub-goals at test time, making it easier to perform multi-step tasks. However, when the cloth configuration varies from that in demonstration sub-goal images, the flow estimation can be disturbed by the variation of the cloth configurations. Hence, the flow-based representation fails to describe the change in the cloth states between the current observation and the sub-goal. Ganapathi *et al.* [18] learns dense visual correspondences between cloths of different size, color, and pose for sequential multi-step cloth manipulation. However, it relies on demonstration observations and actions provided by a human expert at test time. In contrast, our approach only needs general demonstration observations and can perform similar tasks on different cloth configurations without accessing the expert actions.

B. Space-time Attention

The large-scale adoption of the attention mechanism was first proposed by Transformer architecture [19] in natural language processing and achieved great success. The attention mechanism can capture both local and global contexts with less domain-specific inductive bias, making it easily applied to some other domains. Vision Transformer [20] extends Transformer architecture to image classification tasks by splitting images into sequences of image patches. Bertasius *et al.* [21] propose TimeSformer and adapt the standard Transformer to video with space-time attention mechanism. The space-time attention mechanism extracts a video's temporal and spatial features separately with a great speed-accuracy trade-off. By leveraging space-time attention, TimeSformer is also faster to train than traditional CNN video architecture and achieves higher efficiency. By utilizing the space-time attention mechanism, we extract the spatiotemporal features of a sequence of subgoals to capture the instruction information behind the demonstration.

III. APPROACH

A. Problem Formulation

Our goal is to enable a robot to perform sequential multi-step cloth manipulation tasks under the instruction of a general demonstration. Let each task be defined by a sequence of sub-goal observations $\mathcal{G} : \{x_1^g, x_2^g, \dots, x_N^g\}$, these sub-goal observations are obtained from an expert demonstration. Rather than a single goal, we use sub-goals because the cloth of the goal state in sequential multi-step cloth manipulation tasks can be highly self-occluded. A single goal loses part of the information about the full goal state. Unlike prior works [11],

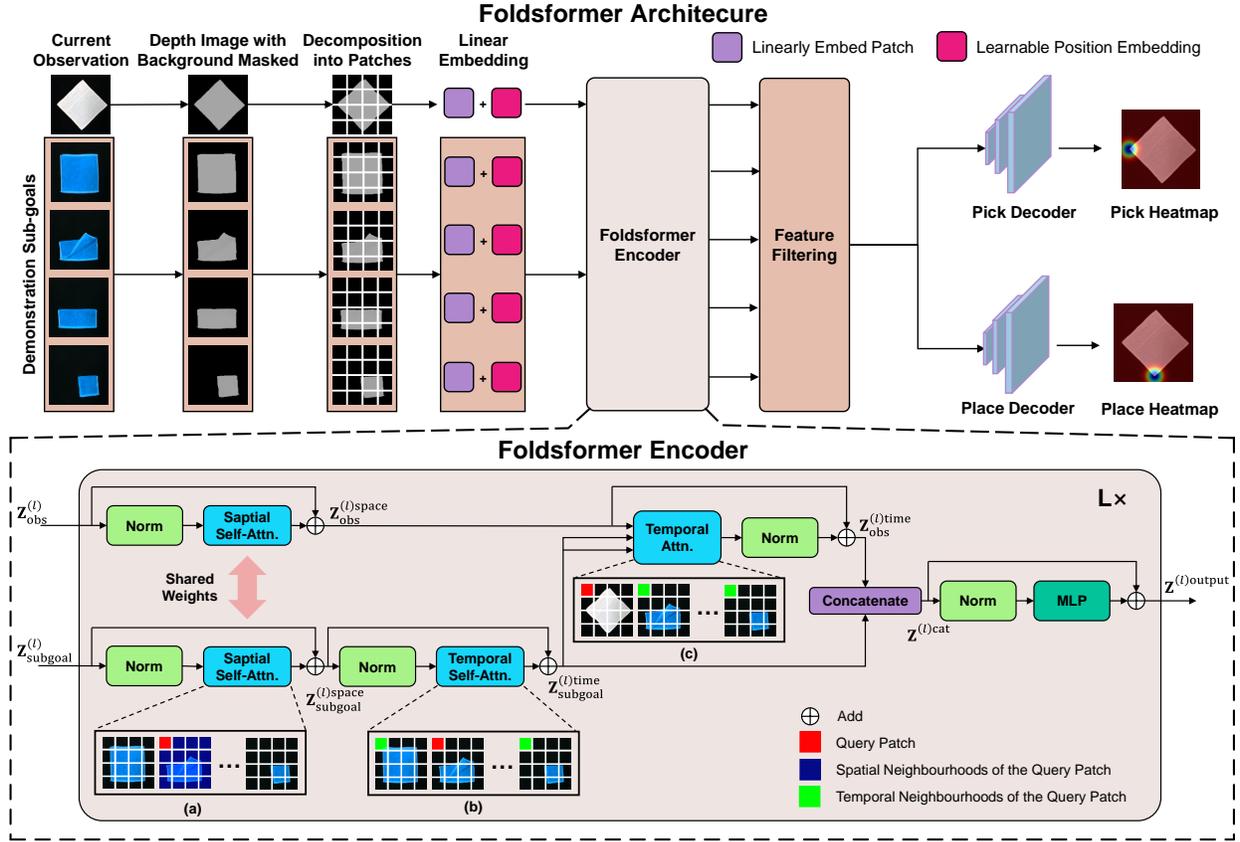


Fig. 2. **Foldsformer Architecture.** Foldsformer takes the current observation depth image and the sub-goal depth image sequence as input and outputs a pick heatmap and a place heatmap that simultaneously estimates a pick point and a place point. In the Foldsformer encoder block, we split the input into two vectors and process them in different ways, and we finally make a fusion to aggregate information. For illustration, we denote the query patch in red, its spatial neighborhoods in blue, and its temporal neighborhoods in green. Patches without color are not used for the attention computation of the red patch. (a) Visualization of the space multi-head self-attention. The self-attention is computed for every patch in all images. (b) Visualization of the time multi-head self-attention. The self-attention is computed for every patch in sub-goal images. (c) Visualization of the time multi-head attention. The self-attention is computed for every patch in the current observation image.

[22], [23], we do not require the cloth configuration in sub-goals to be the same as the cloth configuration at test time. A general demonstration should be applied in different settings.

We formulate the problem as learning a policy π that infers a pick-and-place action $a_t \in \mathcal{A}$ from the current visual observation of the cloth $o_t \in \mathcal{O}$ and the sub-goals \mathcal{G} :

$$\pi(o_t, \mathcal{G}) \rightarrow a_t = (p_{\text{pick}}, p_{\text{place}}) \in \mathcal{A} \quad (1)$$

where p_{pick} is the pixel pick point of the end effector when grasping part of the cloth, and p_{place} is the pixel place point of the end effector when releasing the grasp.

B. Foldsformer

Previous approaches input the current observation o_t and each sub-goal x_i^g sequentially to their goal-conditioned policies [11], [22], [23]. In this way, these policies can not access a complete picture of the multi-step task, resulting in the loss of the temporal information of the sub-goal sequence. Hence, these policies can only capture limited instruction information from the demonstration sub-goals, as they can only access each sub-goal separately. Our insight is that we can allow the policy to access all the sub-goals together and take advantage of the temporal features of all sub-goal images to obtain more instruction information about a multi-step task.

The schematic of the proposed Foldsformer architecture can be found in Fig. 2. Foldsformer is a variant of TimeS-former [21] in sequential multi-step cloth manipulation. Foldsformer takes the current observation o_t and the demonstration sub-goal observation sequence \mathcal{G} as input and outputs a pick heatmap and a place heatmap estimating a pick point and a place point simultaneously. Further details of Foldsformer are described below.

Foldsformer input. Foldsformer takes as input a sequence of depth images $X = \{o_t, \mathcal{G}\} \in \mathbb{R}^{H \times W \times 1 \times (F+1)}$ that consists of the current observation depth image o_t of size $H \times W$ and F sub-goal depth images \mathcal{G} of size $H \times W$. We use depth rather than RGB inspired by [11], because the gap between simulated depth images and real depth images is small. By using depth-only input, our approach can be easily transferred to the real world without any domain randomization.

Decomposition into patches. We decompose each image into N non-overlapping patches, each of size $P \times P$, then we flatten these patches into vectors $\mathbf{x}_{(p,t)} \in \mathbb{R}^{P^2}$ with $p = 1, \dots, N$ denoting spatial locations and $t = 1, \dots, F+1$ denoting an index over images.

Linear Embedding. We map each patch $\mathbf{x}_{(p,t)}$ into an embedding vector $\mathbf{z}_{(p,t)} \in \mathbb{R}^D$ of D dimensions linearly:

$$\mathbf{z}_{(p,t)}^{(0)} = E\mathbf{x}_{(p,t)} + \mathbf{e}_{(p,t)}^{\text{pos}} \quad (2)$$

where $E \in \mathbb{R}^{D \times P^2}$ is a learnable matrix, $\mathbf{e}_{(p,t)}^{pos} \in \mathbb{R}^D$ denotes a learnable position embedding that encodes the spatial and temporal position of each patch.

Foldsformer Encoder. We design our Foldsformer encoder block (see Fig. 2) based on the divided space-time attention mechanism in TimeFormer [21]. Compared to the TimeFormer encoder block, Foldsformer encoder block first splits the input into vectors corresponding to the current observation and vectors corresponding to the demonstration sub-goal sequence, then processes them in different ways and finally makes a fusion of them to aggregate information. We process them separately because the cloth's configuration in the current observation may differ from that in the demonstration sub-goal sequence.

Let $\mathbf{z}^{(l)\text{input}} \in \mathbb{R}^{\frac{H}{P} \times \frac{W}{P} \times D \times (F+1)}$ denote the input vectors of the encoder block l , with $l = 1, \dots, L$ denoting an index over blocks. $\mathbf{z}^{(l)\text{input}}$ can be split in time dimension into vectors $\mathbf{z}_{\text{obs}}^{(l)} \in \mathbb{R}^{\frac{H}{P} \times \frac{W}{P} \times D \times 1}$ that correspond to the current observation image and $\mathbf{z}_{\text{subgoal}}^{(l)} \in \mathbb{R}^{\frac{H}{P} \times \frac{W}{P} \times D \times F}$ that correspond to the sub-goal image sequence.

We firstly compute multi-head self-attention [20] of vectors $\mathbf{z}_{\text{obs}}^{(l)}$ and $\mathbf{z}_{\text{subgoal}}^{(l)}$ with the same weights in space dimension. Each patch and its spatial neighbors within the same image are used for the space self-attention computation (see Fig. 2(a)). We use residual connection after space self-attention to aggregate information:

$$\begin{aligned} \mathbf{z}_{\text{obs}}^{(l)\text{space}} &= \text{SpaceMSA}(\text{LN}(\mathbf{z}_{\text{obs}}^{(l)})) + \mathbf{z}_{\text{obs}}^{(l)} \\ \mathbf{z}_{\text{subgoal}}^{(l)\text{space}} &= \text{SpaceMSA}(\text{LN}(\mathbf{z}_{\text{subgoal}}^{(l)})) + \mathbf{z}_{\text{subgoal}}^{(l)} \end{aligned} \quad (3)$$

where $\text{LN}(\cdot)$ denotes LayerNorm [24], and $\text{SpaceMSA}(\cdot)$ denotes space multi-head self-attention. This space multi-head self-attention extracts spatial information corresponding to the cloth state in the current observation image and each sub-goal image.

We then compute multi-head self-attention of vectors $\mathbf{z}_{\text{subgoal}}^{(l)\text{space}}$ in time dimension. Each patch and its temporal neighbors at the same spatial location are used for the time self-attention computation (see Fig. 2(b)). We use residual connection after time self-attention to aggregate information:

$$\mathbf{z}_{\text{subgoal}}^{(l)\text{time}} = \text{TimeMSA}(\text{LN}(\mathbf{z}_{\text{subgoal}}^{(l)\text{space}})) + \mathbf{z}_{\text{subgoal}}^{(l)\text{space}} \quad (4)$$

where $\text{TimeMSA}(\cdot)$ denotes time multi-head self-attention. This time multi-head self-attention extracts temporal information corresponding to the cloth state changes among all sub-goal images.

After space multi-head self-attention and time multi-head self-attention, we then add a time multi-head attention (see Fig. 2(c)) to fuse the extracted information of the current observation and the sub-goals. The time multi-head attention is computed in time dimension (see Fig. 2(c)), where $\mathbf{z}_{\text{obs}}^{(l)\text{space}}$ serves as the query, $\mathbf{z}_{\text{subgoal}}^{(l)\text{time}}$ serves as the key and the value. We use residual connection after the time multi-head attention to aggregate information:

$$\begin{aligned} \mathbf{z}_{\text{obs}}^{(l)\text{time}} &= \text{LN}(\text{TimeMA}(\mathbf{z}_{\text{obs}}^{(l)\text{space}}, \mathbf{z}_{\text{subgoal}}^{(l)\text{time}} \\ &\quad \mathbf{z}_{\text{subgoal}}^{(l)\text{time}})) + \mathbf{z}_{\text{obs}}^{(l)\text{space}} \end{aligned} \quad (5)$$

where $\text{TimeMA}(\cdot)$ denotes time multi-head attention.

Then, the concatenation $\mathbf{z}^{(l)\text{cat}}$ of $\mathbf{z}_{\text{obs}}^{(l)\text{time}}$ and $\mathbf{z}_{\text{subgoal}}^{(l)\text{time}}$ is projected and passed through an MLP, using residual connection:

$$\begin{aligned} \mathbf{z}^{(l)\text{cat}} &= \text{Concatenate}(\mathbf{z}_{\text{obs}}^{(l)\text{time}}, \mathbf{z}_{\text{subgoal}}^{(l)\text{time}}) \\ \mathbf{z}^{(l)\text{output}} &= \text{MLP}(\text{LN}(\mathbf{z}^{(l)\text{cat}})) + \mathbf{z}^{(l)\text{cat}} \end{aligned} \quad (6)$$

The output of the Foldsformer encoder block $\mathbf{z}^{(l)\text{output}}$ keeps the same dimension as the input of the Foldsformer encoder block $\mathbf{z}^{(l)\text{input}}$, so these encoder blocks can be stacked in a sequence to make the network deep. Our encoder consists of several Foldsformer encoder blocks.

Feature filtering. In the Foldsformer encoder output, we only use the vectors that correspond to the current observation and vectors that correspond to the last image in the sub-goal image sequence, i.e., the goal image. Let $\mathbf{z}_{(t)}^{\text{output}} \in \mathbb{R}^{\frac{H}{P} \times \frac{W}{P} \times D \times (F+1)}$ denote the encoder output, where $t = 1, \dots, F+1$ denotes an index over images. We choose $\mathbf{z}_{(1)}^{\text{output}}$ and $\mathbf{z}_{(F+1)}^{\text{output}}$ and concatenate them for the decoder input:

$$\mathbf{z}^{\text{filter}} = \text{Concatenate}(\mathbf{z}_{(1)}^{\text{output}}, \mathbf{z}_{(F+1)}^{\text{output}}) \quad (7)$$

In this way, the decoder only takes as input the feature of the current observation and the goal image to infer an action, forming closed loop feedback.

Decoder. We build a pick decoder and a place decoder to estimate a pick and place point, respectively. The two decoders have the same architecture. Similar to [25], we consider a progressive upsampling strategy that alternates convolutional layers and upsampling operations. The two decoders take $\mathbf{z}^{\text{filter}}$ as input, then output a predicted pick heatmap Q_{pick} and a predicted place heatmap Q_{place} that have the same size as the input images. From the two heatmaps, we can get the optimal pick point and the optimal place point:

$$\begin{aligned} p_{\text{pick}} &= \text{argmax}_p Q_{\text{pick}}(p) \\ p_{\text{place}} &= \text{argmax}_p Q_{\text{place}}(p) \end{aligned} \quad (8)$$

Loss. We supervise the whole network using the binary cross-entropy(BCE) loss between predicted heatmaps Q_{pick} , Q_{place} and ground truth heatmaps $Q_{\text{pick}}^{\text{gt}}$, $Q_{\text{place}}^{\text{gt}}$, and compute the sum of them:

$$\mathcal{L} = \text{BCE}(Q_{\text{pick}}, Q_{\text{pick}}^{\text{gt}}) + \text{BCE}(Q_{\text{place}}, Q_{\text{place}}^{\text{gt}}) \quad (9)$$

C. Implementation Details

We collect training data and evaluate Foldsformer and the baselines in the SoftGym suite [26], a particle-based simulator built on Nvidia Flex. In our simulation environment, a top-down camera provides depth and RGB images of size 224×224 , and the cloth is put at the center of the camera view.

Our training data consists of samples generated by random actions and expert demonstrations. For collecting random training samples, we first generate 1000 random initial configurations of rectangular cloths whose size and pose are randomized. Then we collect 6000 trajectories of length 8 using random actions. Inspired by [11], 80% of the random actions are biased to pick corners, and 20% of the pick actions

TABLE I. Mean particle distance error (mm) and inference time (sec) on four sequential multi-step cloth manipulation tasks in simulation. The best results are bolded. We bold the results of Foldsformer (Full Approach) and Foldsformer (NoTimeAttn) in the *Double Straight* task that have a significant overlap.

Approach	<i>Double Triangle</i>	<i>Double Straight</i>	<i>All Corners Inward</i>	<i>Corners Edges Inward</i>	Inference Time
FabricFlowNet [11]	102.20 ± 19.39	85.34 ± 19.67	33.64 ± 11.76	36.95 ± 9.95	~ 0.002
Lee <i>et al.</i> [10]	109.82 ± 39.96	114.71 ± 26.06	27.21 ± 11.47	70.67 ± 22.24	~ 0.061
Fabric-VSF [12]	114.62 ± 35.46	116.94 ± 24.52	46.05 ± 23.38	51.82 ± 16.65	~ 6.512
Foldsformer (Goal-Conditioned)	93.30 ± 43.35	80.81 ± 38.39	49.46 ± 19.55	56.66 ± 43.25	~ 0.007
Foldsformer (NoTimeAttn)	59.86 ± 44.58	60.90 ± 47.52	9.11 ± 3.09	22.97 ± 20.42	~ 0.005
Foldsformer (Full Approach)	19.64 ± 17.08	59.09 ± 44.82	3.06 ± 1.83	8.11 ± 7.96	~ 0.011

are uniformly sampled over the cloth. We use corner bias because most sequential multi-step cloth manipulation tasks tend to pick corners. For collecting expert demonstrations, we use scripted demonstrator policies to get 100 trajectories of expert demonstrations per task. In SoftGym, the cloth is modeled as a grid of particles, and the ground truth position of each particle can be accessed. Thus, we can easily get the true corners and script a demonstrator.

In each trajectory, we take a sequence of actions and record the initial observation o_0 , actions $a_i = (p_{\text{pick}}^{(i)}, p_{\text{place}}^{(i)})$, and the next observation o_i after action a_i , with $i = 1, 2, \dots, M$ denoting an index over actions. Mathematically, a trajectory \mathcal{T} of length M consists of a sequence of observations and actions:

$$\mathcal{T} = \{o_0, a_1, o_1, \dots, o_{M-1}, a_M, o_M\} \quad (10)$$

To train Foldsformer faster, we split each trajectory whose length is M into sub-trajectories whose length is K . Let $\mathcal{S}_i = \{o_i, a_{i+1}, o_{i+1}, \dots, o_{i+K-1}, a_{i+K}, o_{i+K} | i = 0, 1, \dots, M - K\}$ denote a sub-trajectory. From a sub-trajectory, we can get K training samples: the input $(o_{i+j}, o_i, o_{i+1}, \dots, o_{i+K})$ and the action a_{i+j+1} used to supervise the network with $j = 0, 1, \dots, K - 1$. In this way, we can get $(M - K + 1)H$ training samples from each trajectory of length M , making our approach sample-efficient.

In this work, the patch size (P), the embedding vectors’ dimension (D), the number of encoder blocks (L), and the length of the sub-trajectory (K) are set to 16, 256, 8, 4, respectively. Each encoder alternates 5 convolutional layers and 4 bilinear upsampling layers. The kernel sizes, strides, and filter sizes of the convolutional layers are [1, 1, 1, 1, 1], [1, 1, 1, 1, 1], [256, 256, 128, 128, 1], respectively. The scale factors of the bilinear upsampling layers are set to 2. We train Foldsformer with a batch size of 32, a learning rate of 0.0001, and an Adam optimizer [27] for 20 epochs. It takes about 9 hours on an NVIDIA RTX 3090 GPU, and the BCE loss drops to about 0.02.

IV. EXPERIMENTS

A. Simulation Setup

As mentioned in Sec. III-C, we use SoftGym [26] simulator for training and evaluation. We evaluate Foldsformer and the baselines on 4 sequential multi-step cloth manipulation tasks, which can be achieved by 2, 3, 4, and 4 actions, respectively (see Fig. 3 for details.): (i) *Double Triangle*; (ii) *Double Straight*; (iii) *All Corners Inward*; (iv) *Corners Edges Inward*. We consider 40 different cloth configurations per task, combining 10 different sizes and 4 different rotation

angles. For task (i), task (ii), and task (iv), we use square cloth, the 10 different sizes range from 31.25cm × 31.25cm to 36.875cm × 36.875cm divided into 0.625cm × 0.625cm intervals. For task (ii), we use rectangular cloth, and the 10 different sizes range from 31.25cm × 27.5cm to 36.875cm × 32.5cm divided into 0.625cm × 0.5cm intervals. The 4 different rotation angles are 0°, 30°, 45° and 60°. Given a sequence of demonstration sub-goal images where the square cloth size is 34.375cm × 34.375cm, the rectangle cloth size is 34.375cm × 27.5cm and the rotation angle is 0°, a policy needs to perform on all 40 different cloth configurations. Our error metric is the mean particle position error between the cloth states achieved by the policy and a scripted demonstrator mentioned in Sec. III-C that can access the ground truth particle positions.

We compare our approach Foldsformer to three previous state-of-the-art baselines for sequential multi-step cloth manipulation: FabricFlowNet [11], which learns a flow-based policy; Lee *et al.* [10], which uses offline RL for arbitrary-goal cloth folding; Fabric-VSF [12], which learns a forward cloth dynamic model in pixel space and use it for planning. We adapt FabricFlowNet to a single-arm variant. Although the original approaches [10], [12] use only a single goal image, we provide sub-goals for them.

To further evaluate the importance of the components of our approach, we conduct two ablations in simulation: A standard goal-conditioned variant of our approach (“Goal-Conditioned”), which only takes as input the current observation and each sub-goal separately; An variant of our approach (“NoTimeAttn”), which removes time attention in the Foldsformer encoder block described in Sec. III-B.

For each task, we fine-tune our approach, the baselines, and the ablations with 100 expert demonstrations.

B. Simulation Results

Simulation results can be found in Table. I. Foldsformer achieves the lowest mean particle error over all tasks and shows a larger improvement over the baselines. Fig. 3 shows representative qualitative results on 5 of all 40 different clothes. As can be seen, even when the size and pose of the cloth are different from that of the cloth in demonstration sub-goals, Foldsformer performs all tasks successfully, while the baselines fail to complete most of the tasks. Although the inference time¹ of Foldsformer is slower than FabricFlowNet (~0.002s vs. ~0.011s), it is still acceptable in real-time applications.

¹The inference time is tested on one machine with Intel i9-9900K CPU (3.60GHz) and NVIDIA GeForce RTX 2080Ti



Fig. 3. **Qualitative Results in Simulation Experiments.** We evaluate Foldsformer and the baselines on 4 sequential multi-step cloth manipulation tasks: (a) *Double Triangle*. (b) *Double Straight*. (c) *All Corners Inward*. (d) *Corners Edges Inward*.

For the ablations, the “Goal-Conditioned” ablation validates the significance of allowing a policy to access all the sub-goals together. A policy can capture more instruction information by accessing all the sub-goals together rather than accessing them separately. The “NoTimeAttn” ablation demonstrates the benefits of time attention to leverage the temporal features of the sub-goal images. Previous approaches input the sub-goals separately, leading to the loss of these temporal features.

C. Real World Setup

Fig. 4(a) shows the robot system in the real world. We use a 7-DOF Franka Emika Panda robot arm with a standard two-finger panda gripper. The pick and place trajectories are planned by using MoveIt! [28]. We obtain RGB-D images from an Intel Realsense D435i camera mounted at the end of the arm. Before capturing images, the arm moves to a fixed pose where the camera can capture an effective workspace of $56\text{cm} \times 56\text{cm}$. To narrow the gap between simulated depth images and real-world depth images, we crop the center of the raw depth images from the camera and segment the cloth by color-thresholding the background. Similar to [11], we also add the difference between the simulated surface and the mean depth of the support surface in the real world to align the depth between real and simulated depth images. Through these simple techniques, our approach trained in simulation can be transferred to the real world successfully without additional training or domain randomization.

We use 3 square cloths of different size: a $20\text{cm} \times 20\text{cm}$ cloth, a $30\text{cm} \times 30\text{cm}$ cloth and a $35\text{cm} \times 35\text{cm}$ cloth. We vary the pose by placing the cloths in 3 rotation angles (0° , 30° , and 45°). We also use a T-shirt ($50\text{cm} \times 35\text{cm}$) and a pair of shorts ($40\text{cm} \times 32\text{cm}$) which are unseen during training. All clothes can be found in Fig. 4. We complete the same four tasks in simulation using the 3 square cloths. In each task, given a sequence of demonstration images where the cloth’s size is $30\text{cm} \times 30\text{cm}$, and the rotation angle is 0° , Foldsformer should perform the task on all 9 combinations of

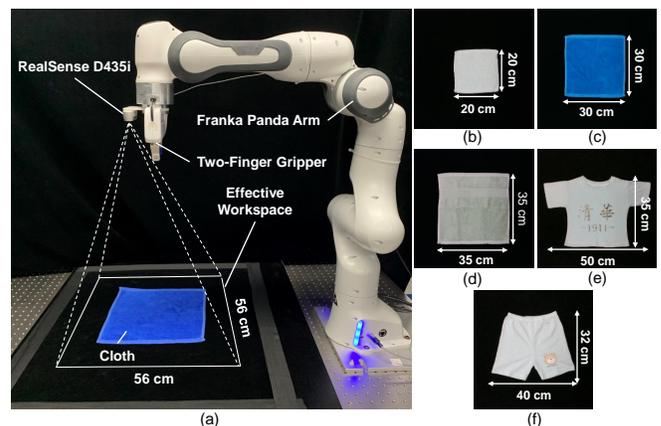


Fig. 4. **Experimental Setup.** (a) Robot system. (b) A $20\text{cm} \times 20\text{cm}$ cloth. (c) A $30\text{cm} \times 30\text{cm}$ cloth. (d) A $35\text{cm} \times 35\text{cm}$ cloth. (e) A $50\text{cm} \times 35\text{cm}$ T-shirt. (f) A pair of shorts ($40\text{cm} \times 32\text{cm}$).

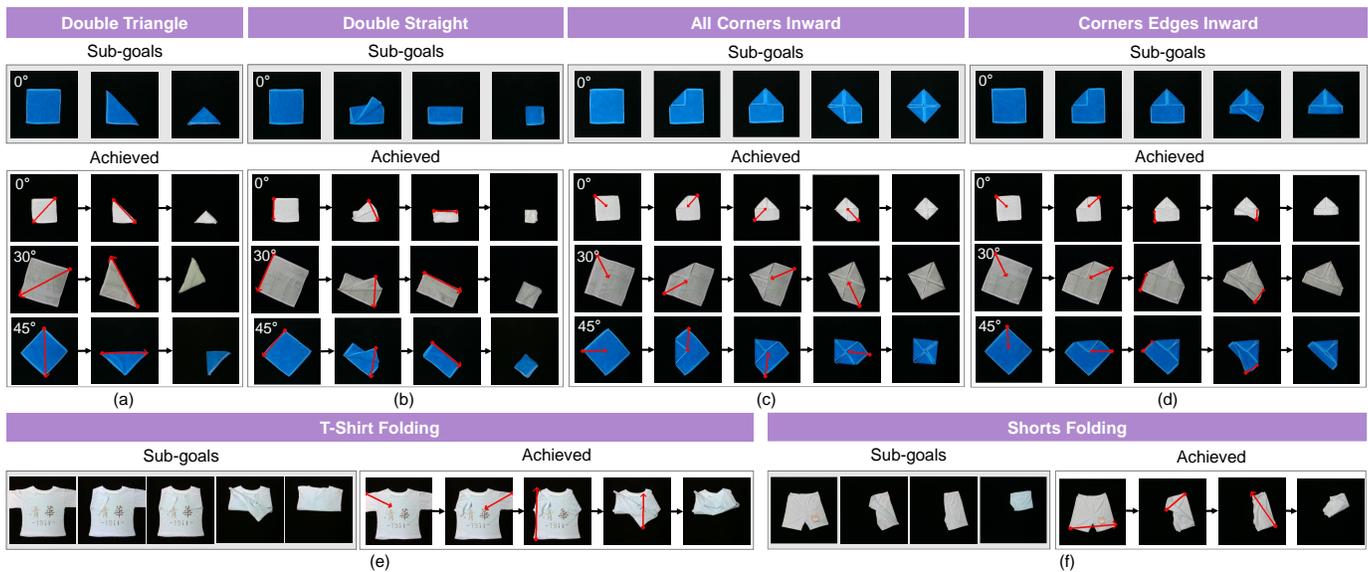


Fig. 5. **Qualitative Results in Real World Experiments.** Foldsformer performs 4 sequential multi-step cloth manipulation tasks on square cloths in the real world: (a) *Double Triangle*. (b) *Double Straight*. (c) *All Corners Inward*. (d) *Corners Edges Inward*. Foldsformer also completes a T-shirt folding task and a shorts folding task where the T-shirt and shorts are unseen cloth shapes during training: (e) *T-shirt Folding*. (f) *Shorts Folding*. Pick and place actions are visualized as red arrows.

different cloth size and pose. We also test the generalization of Foldsformer by completing a T-shirt folding task and a shorts folding task where the T-shirt and the shorts (see Fig. 5) are unseen cloth shapes during training. Similar to previous works [10], [11], [29], we evaluate the performance quantitatively by using Mean Intersection over Union (MIoU) between the cloth masks achieved by Foldsformer and a human demonstrator and a penalty for wrinkles (WR) that calculates the fraction of pixels inside the cloth mask detected as edges by the Canny edge detector [30].

D. Real World Results

Table. II shows the quantitative results, and Fig. 5 shows the part of our representative qualitative results. As can be seen, Foldsformer performs all the tasks successfully, even when the cloth’s size and pose differ from the cloth in demonstration sub-goals. Foldsformer can also generalize to the T-shirt and shorts, which are unseen during training. By accessing the complete picture of the task that contains all sub-goals together, Foldsformer can learn the instruction information and be robust to the variation of cloth configurations. Our encoder enables the instruction information learned from a general demonstration can guide Foldsformer perform similar tasks on different cloth configurations. Hence, only a general demonstration is required by human demonstration and can

be applied to different settings, which is more practical and efficient than previous approaches. Besides, since the cloth is often self-occluded and partially observable, providing all sub-goals together is helpful for an agent to obtain the global knowledge of a multi-step task and then infer proper chained actions. Videos of robot executions can be found on our project website.

V. LIMITATIONS AND FUTURE WORK

Failure cases in our experiments mainly come from grasping failures. Since we use a pick-and-place action primitive in 2D pixel space, it’s hard to determine the number of the grasped cloth layers. This grasp quality is crucial in the real-world setting, especially for *Double Straight* and *Corners Edges Inward* tasks. Fig. 6 illustrates this failure case. The cloth size should also be much smaller than the robot arm’s workspace due to physical reachability limitations. Meanwhile, the sim-to-real transfer relies on color-thresholding the background, which is not feasible in some cluttered backgrounds. Further, although Foldsformer can successfully perform tasks on clothes of different configurations with a single general demonstration, the demonstration sub-goals still need to be generated by expert demonstrations. In future work, we will try

TABLE II. MIoU and mean WR on six sequential multi-step cloth manipulation tasks in the real world

Task	MIoU	mean WR
<i>Double Triangle</i>	0.885	0.022
<i>Double Straight</i>	0.872	0.035
<i>All Corners Inward</i>	0.911	0.042
<i>Corners Edges Inward</i>	0.923	0.044
<i>T-shirt Folding</i>	0.801	0.106
<i>Shorts Folding</i>	0.676	0.061

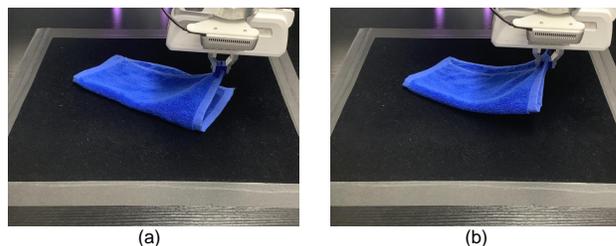


Fig. 6. **Failure Case.** Failure cases in our experiments mainly come from grasping failures. (a) Desired Grasp. The gripper is supposed to grasp one layer of the towel. (b) Executed Grasp. The gripper grasps two layers of the towel.

to make a fusion of visual information and tactile signals [31] to improve the grasp quality and reliability. We will also explore generating the sub-goals automatically.

VI. CONCLUSION

In this paper, we present a novel framework named Foldsformer for sequential multi-step cloth manipulation. We utilize the space-time attention mechanism to capture the instruction information behind a general demonstration that can guide Foldsformer complete similar tasks on different cloth configurations. Experiments results indicate that Foldsformer is effective and practical for sequential multi-step cloth manipulation tasks both in simulation and in the real world. Furthermore, Foldsformer can also generalize to unseen cloth shapes despite only being trained on rectangular cloths.

REFERENCES

- [1] A. Doumanoglou, J. Stria, G. Peleka, I. Mariolis, V. Petrík, A. Kargakos, L. Wagner, V. Hlaváč, T.-K. Kim, and S. Malassiotis, "Folding clothes autonomously: A complete pipeline," *IEEE Transactions on Robotics*, vol. 32, no. 6, pp. 1461–1478, 2016.
- [2] C. Bersch, B. Pitzer, and S. Kammel, "Bimanual robotic cloth manipulation for laundry folding," in *2011 IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2011, pp. 1413–1419.
- [3] B. Thananjeyan, A. Garg, S. Krishnan, C. Chen, L. Miller, and K. Goldberg, "Multilateral surgical pattern cutting in 2d orthotropic gauze with deep reinforcement learning policies for tensioning," in *2017 IEEE International Conference on Robotics and Automation (ICRA)*, 2017, pp. 2371–2378.
- [4] E. Torgerson and F. W. Paul, "Vision-guided robotic fabric manipulation for apparel manufacturing," *IEEE Control Systems Magazine*, vol. 8, no. 1, pp. 14–20, 1988.
- [5] N. Essahbi, B. C. Bouzgarrou, and G. Gogu, "Soft material modeling for robotic manipulation," in *Applied Mechanics and Materials*, vol. 162. Trans Tech Publ, 2012, pp. 184–193.
- [6] D. Seita, P. Florence, J. Tompson, E. Coumans, V. Sindhwani, K. Goldberg, and A. Zeng, "Learning to rearrange deformable cables, fabrics, and bags with goal-conditioned transporter networks," in *2021 IEEE International Conference on Robotics and Automation (ICRA)*, 2021, pp. 4568–4575.
- [7] D. Seita, A. Ganapathi, R. Hoque, M. Hwang, E. Cen, A. K. Tanwani, A. Balakrishna, B. Thananjeyan, J. Ichnowski, N. Jamali, K. Yamane, S. Iba, J. Canny, and K. Goldberg, "Deep imitation learning of sequential fabric smoothing from an algorithmic supervisor," in *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2020, pp. 9651–9658.
- [8] X. Wang, J. Zhao, X. Jiang, and Y.-H. Liu, "Learning-based fabric folding and box wrapping," *IEEE Robotics and Automation Letters*, vol. 7, no. 2, pp. 5703–5710, 2022.
- [9] G. Salhotra, I.-C. A. Liu, M. Dominguez-Kuhne, and G. S. Sukhatme, "Learning deformable object manipulation from expert demonstrations," *IEEE Robotics and Automation Letters*, vol. 7, no. 4, pp. 8775–8782, 2022.
- [10] R. Lee, D. Ward, V. Dasagi, A. Cosgun, J. Leitner, and P. Corke, "Learning arbitrary-goal fabric folding with one hour of real robot experience," in *Conference on Robot Learning*. PMLR, 2021, pp. 2317–2327.
- [11] T. Weng, S. M. Bajracharya, Y. Wang, K. Agrawal, and D. Held, "Fabricflownet: Bimanual cloth manipulation with a flow-based policy," in *Conference on Robot Learning*. PMLR, 2022, pp. 192–202.
- [12] R. Hoque, D. Seita, A. Balakrishna, A. Ganapathi, A. Tanwani, N. Jamali, K. Yamane, S. Iba, and K. Goldberg, "VisuoSpatial Foresight for Multi-Step, Multi-Task Fabric Manipulation," in *Robotics: Science and Systems (RSS)*, 2020.
- [13] Y. Yamakawa, A. Namiki, and M. Ishikawa, "Motion planning for dynamic folding of a cloth with two high-speed robot hands and two high-speed sliders," in *2011 IEEE International Conference on Robotics and Automation*, 2011, pp. 5486–5491.
- [14] S. Miller, J. Van Den Berg, M. Fritz, T. Darrell, K. Goldberg, and P. Abbeel, "A geometric approach to robotic laundry folding," *The International Journal of Robotics Research*, vol. 31, no. 2, pp. 249–267, 2012.
- [15] W. Yan, A. Vangipuram, P. Abbeel, and L. Pinto, "Learning predictive representations for deformable objects using contrastive estimation," in *Conference on Robot Learning*. PMLR, 2021, pp. 564–574.
- [16] X. Ma, D. Hsu, and W. S. Lee, "Learning latent graph dynamics for visual manipulation of deformable objects," in *2022 International Conference on Robotics and Automation (ICRA)*, 2022, pp. 8266–8273.
- [17] X. Lin, Y. Wang, Z. Huang, and D. Held, "Learning visible connectivity dynamics for cloth smoothing," in *Conference on Robot Learning*. PMLR, 2022, pp. 256–266.
- [18] 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 dense visual correspondences in simulation to smooth and fold real fabrics," in *2021 IEEE International Conference on Robotics and Automation (ICRA)*, 2021, pp. 11 515–11 522.
- [19] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," *Advances in neural information processing systems*, vol. 30, 2017.
- [20] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly *et al.*, "An image is worth 16x16 words: Transformers for image recognition at scale," in *International Conference on Learning Representations*, 2020.
- [21] G. Bertasius, H. Wang, and L. Torresani, "Is space-time attention all you need for video understanding?" in *International Conference on Machine Learning*. PMLR, 2021, pp. 813–824.
- [22] A. Nair, D. Chen, P. Agrawal, P. Isola, P. Abbeel, J. Malik, and S. Levine, "Combining self-supervised learning and imitation for vision-based rope manipulation," in *2017 IEEE International Conference on Robotics and Automation (ICRA)*, 2017, pp. 2146–2153.
- [23] D. Pathak*, P. Mahmoudieh*, M. Luo*, P. Agrawal*, D. Chen, F. Shentu, E. Shelhamer, J. Malik, A. A. Efros, and T. Darrell, "Zero-shot visual imitation," in *International Conference on Learning Representations*, 2018. [Online]. Available: <https://openreview.net/forum?id=BKisuzWRW>
- [24] J. L. Ba, J. R. Kiros, and G. E. Hinton, "Layer normalization," *arXiv preprint arXiv:1607.06450*, 2016.
- [25] S. Zheng, J. Lu, H. Zhao, X. Zhu, Z. Luo, Y. Wang, Y. Fu, J. Feng, T. Xiang, P. H. Torr *et al.*, "Rethinking semantic segmentation from a sequence-to-sequence perspective with transformers," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 6881–6890.
- [26] X. Lin, Y. Wang, J. Olkin, and D. Held, "Softgym: Benchmarking deep reinforcement learning for deformable object manipulation," in *Conference on Robot Learning*. PMLR, 2021, pp. 432–448.
- [27] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.
- [28] S. Chitta, I. Sucas, and S. Cousins, "Moveit![ros topics]," *IEEE Robotics & Automation Magazine*, vol. 19, no. 1, pp. 18–19, 2012.
- [29] 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.
- [30] J. Canny, "A computational approach to edge detection," *IEEE Transactions on pattern analysis and machine intelligence*, no. 6, pp. 679–698, 1986.
- [31] S. Tirumala, T. Weng, D. Seita, O. Kroemer, Z. Temel, and D. Held, "Learning to Singulate Layers of Cloth based on Tactile Feedback," in *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2022.