

Appendix

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25

Contents

1	Implementation Details	2
1.1	Entity Encoding	2
1.2	Details on Training Set Feasibility Predictor	2
1.3	Details on the Attention Structure for Planning with Set Representation	2
1.4	Network Architectures	3
1.5	Training Details	4
2	Details on Simulation Experiments	5
2.1	Hyperparameters for Simulation Dough	5
2.2	Hyperparameters for DBSCAN	5
2.3	Hyperparameters for PASTA	5
2.4	Receding Horizon Planning for CRS-Twice	6
3	Details on Real World Experiments	6
3.1	Heuristic Policies	6
3.2	Procedure for Resetting the Dough	7
3.3	Procedure for the Human Baseline	7
3.4	Procedure for Making the Real Dough	7
4	Additional Experiments	7
4.1	Ablation Studies	7
4.2	Visualization of the Latent Space	8
4.3	Runtime of PASTA	9
4.4	Additional metrics for real world experiments	9
4.5	Robustness of PASTA	9
5	Further Discussion on Limitations and Future Work	10

26 1 Implementation Details

27 1.1 Entity Encoding

28 To train our point cloud variational autoencoder [1], we normalize the point cloud of each entity
29 P_i to be centered at the origin, i.e. $\bar{P}_i = P_i - t_i$, where t_i is the mean of all points in P_i . We then
30 encode each centered \bar{P}_i into a latent encoding: $z_i = \phi(\bar{P}_i)$. Our latent representation $u_i = [z_i, t_i]$
31 consists of the encoding of the point cloud’s shape z_i and the position of the center of the point
32 cloud t_i . We model the point cloud’s position t_i explicitly such that the learned latent embedding
33 z_i can focus on shape variation alone and the model that plans over u_i can still reason over point
34 clouds at different spatial locations. During training, we record the 3D bounding box of all training
35 data $t_{min}, t_{max} \in \mathbb{R}^3$, and we sample from $[t_{min}, t_{max}]$ during planning. We denote this combined
36 distribution of $u = [z, t]$ as P_u .

37 1.2 Details on Training Set Feasibility Predictor

38 **Hard Negative Samples** Suppose skill k takes N_k latent vectors from observation \hat{U}^o and M_k latent
39 vectors from goal \hat{U}^g as input. To generate random pairs of observations and goals as negative
40 samples for the feasibility predictor, we can sample each latent point cloud representation u_i by
41 sampling the shape z_i from the VAE prior p_z and sampling the position t_i from the distribution of
42 positions in the training dataset. Such random negative samples are used similarly in DiffSkill [2].
43 However, as the combined dimension of the set representation becomes larger compared to a flat
44 representation, we need a way to generate harder negative samples. To do so, for a positive pair of
45 set representation $(\{u_i^o\}, \{u_j^g\})$, we randomly replace one of the entities u_i^o or u_j^g with a random
46 sample in the latent space and use it as a negative sample. Our ablation results show that this way of
47 generating hard negative samples is crucial for training our set feasibility predictor.

48 **Noise on Latent Vectors** During training of the feasibility predictor, for each of the input latent
49 vector $u = [z, t]$, where $z \in \mathbb{R}^{D_z}$ is the latent encoding of the shape and $t \in \mathbb{R}^3$ is the 3D position
50 of the point cloud, we add a Gaussian noise to each part, i.e. $\hat{z} = z + \sigma_z \epsilon$ and $\hat{t} = t + \sigma_t \epsilon$, where
51 $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. The amount of noise determines the smoothness of the feasibility landscape. Without
52 any noise, planning with gradient descent with the feasibility function becomes much harder.

53 1.3 Details on the Attention Structure for Planning with Set Representation

54 Given the initial latent set observation U^{obs} with N_o components, and the skill sequences k_1, \dots, k_H ,
55 in this section we describe how to generate the attention structure. We denote the latent set represen-
56 tation at step h as U^h , $h = 1 \dots H$ and define $U^0 = U^{obs}$. As skill k_h takes in N_{k_h} components as
57 input and M_{k_h} components as output, by calculation we know that U^h has $N_o + \sum_{i=1}^h M_{k_i} - N_{k_i}$
58 components. From now on, we denote $|U^h| = N_h$ and $U^h = \{u_{h,1}, \dots, u_{h,N_h}\}$. As the skill k_h
59 only applies to a subset of the input U^{h-1} , we now formally define the attention structure at step h
60 to be I^h , which consists of a list of indices, each of length N_{k_h} , such that I^h selects a subset from
61 U^{h-1} to be the input to the feasibility predictor, i.e.

$$\hat{U}^{h-1} = U_{I^h}^{h-1} \subseteq U^{h-1}.$$

62 However, enumerating all I^h is infeasible, as there are $C_{N_h}^{N_{k_h}}$ combinations for each step. Fortunately,
63 we do not have to enumerate all different structures. The insight here is that, for each attention
64 structure I_1, \dots, I_H , we will perform a low-level optimization. In this low-level optimization, we will
65 first initialize all the latent vectors to be optimized from P_u and then perform gradient descent on
66 them. As many of the attention structures yield topologically equivalent tree structures (An example
67 of such a tree is illustrated in Figure 2c of the main paper), and each latent vector in the tree is
68 sampled independently from the same distribution P_u , these topologically equivalent tree structures
69 result in the same optimization process. As such, we do not need to exhaust all of such attention
70 structures.

71 Instead of enumerating each topologically different structure and then sampling multiple initializations
 72 for the low-level optimization, we randomly sample sequences (I_1, \dots, I_H) and perform low-level
 73 gradient-descent optimization on all the samples. In this way, with enough samples, we will be able
 74 to cover all attention structures.

75 Now, we can sequentially build up the subgoals latent set representation during planning. Specifically,
 76 assuming that we have constructed the previous latent set representation U^{h-1} , we will now describe
 77 the procedure for constructing U^h , as well as the predicted feasibility for the current skill k_h ,
 78 i.e. $f_{k_h}(\hat{U}^{o,h}, \hat{U}^{g,h})$, where $\hat{U}^{o,h}, \hat{U}^{g,h}$ are the subset of U^{h-1} and U^h attended by the feasibility
 79 predictor. First, we generate I_h by randomly choosing the index of a subset from U_{h-1} . U^h are
 80 composed of two parts: The first part is the latent vectors generated by applying the skill k_h . For
 81 this part, we will create a set of new vectors $u_{h,0}, \dots, u_{h,M_{h_k}}$. This part of the latent vectors will be
 82 attended by the feasibility predictor as $\hat{U}^{g,h} = \{u_{h,0}, \dots, u_{h,M_{h_k}}\}$. The second part of U^h comes
 83 from the previous latent set vectors that are not modified by the skill, i.e. $U^{h-1} \setminus U_{I_h}^{h-1}$, and U^h is
 84 the addition of both parts, i.e.

$$U^h = \hat{U}^{g,h} \cup (U^{h-1} \setminus U_{I_h}^{h-1})$$

85 In this way, we can sequentially build up U^h from U^{h-1} , and U^0 is simply U^{obs} . At the same time,
 86 we have determined our attention structure and the feasibility prediction. Our objective can thus be
 87 written as

$$\arg \min_{\mathbf{k}, \mathbf{I}, \mathbf{U}} J(\mathbf{k}, \mathbf{I}, \mathbf{U}) = \prod_{h=1}^H f_{k_h}(\hat{U}^{o,h}, \hat{U}^{g,h}) \exp(-C(U^H, U^g)), \quad (1)$$

88 where \mathbf{U} is the set union of all latent vectors to be optimized.

89 1.4 Network Architectures

90 **Set Feasibility Predictor** We use a Multi-Layer Perceptron (MLP) with ReLU activations for our
 91 feasibility predictor. We apply max-pooling to the transformed latent vectors of \hat{U}^o and \hat{U}^g to achieve
 permutation invariance. Below is our architecture:

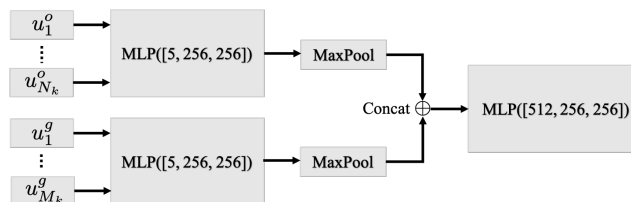


Figure 1: Architecture for the set feasibility predictor

92

93 **Set Cost Predictor** We use a 3-layer MLP with a hidden dimension of 1024 and ReLU activations.

94 **Set Policy** The set point cloud policy for the k^{th} skill π_k takes in an observed point cloud P^{obs} , a
 95 goal point cloud P^{goal} , and a tool point cloud P_k^{tool} and outputs an action at each timestep to control
 96 the tool directly. The tool point cloud P^{tool} is obtained by sampling points on the mesh surface of the
 97 tool and then transforming these points to the same camera frame as the P^{obs} and P^{goal} , assuming
 98 the pose of the tool is known from the robot state. Instead of taking latent vectors as input, the policy
 99 functions directly in point cloud space, which allows it to handle times when spatial abstraction is
 100 ambiguous. For instance, during cutting and merging, the number of dough components gradually
 101 increases or decrease, during which the latent set representation is not changing smoothly while
 102 the point clouds change smoothly during the process. We later show the advantage of using point
 103 clouds directly as the policy input. We concatenate each point's (x, y, z) coordinates with a one-hot
 104 encoding to indicate whether the point belongs to the observation, tool, or goal, and we input the
 105 points into a PointNet++ [3] encoder followed by an MLP which outputs the action. We use a point

	LiftSpread	GatherMove	CutRearrange	CRS + CRS-Twice
# of initial configurations	200	200	1500	1200
# of target configurations	200	200	1500	1200
# of training trajectories	1800	1800	1350	1080
# of testing trajectories	200	200	150	120
# of total trajectories	2000	2000	1500	1200
# of total transitions	1e5	1e5	7.5e4	6e4

Table 1: Summary of training/testing data

cloud for the tool to allow the PointNet encoder to reason about the interaction between the tool and the dough in the same space. We use PyTorch Geometric’s [4] implementation of PointNet++ and with the following list of modules in our encoder.

SAModule(0.5, 0.05, MLP([3+3, 64, 64, 128]))
 SAModule(0.25, 0.1, MLP([128+3, 128, 128, 256]))
 GlobalSAModule(MLP([256+3, 256, 128, 512, 1024]))

The MLP following the encoder consists of hidden dimensions [1024, 512, 256] and ReLU activations.

1.5 Training Details

Training data. We inherit the data generation procedure from DiffSkill [2]: first, we randomly generate initial and target configurations. The variations in these configurations include the location, shape, and size of the dough and the location of the tool. We then sample a specific initial configuration and a target configuration and perform gradient-based trajectory optimization to obtain demonstration data. For each task, the demonstration data consists of all the transitions from executing the actions outputted by the trajectory optimizer. We perform a train/test split on the dataset and select 5 configurations in the test split for evaluating the performance for all the methods. More information about training and testing data can be found in Table 1.

Point cloud VAE. We train our point cloud VAE by maximizing the evidence lower bound (ELBO). For a dataset of observations $P(X)$, which consists of the segmented point cloud of each entity in the scene, we optimize the following objective:

$$\mathcal{L}_{VAE} = \mathbb{E}_{Q_\phi(z|x)} [\log P_\psi(X|z)] - D_{KL}(Q_\phi(z|X)||p(z)) \quad (2)$$

where $Q_\phi(z|X)$ is the encoder modeled as a diagonal Gaussian, $P_\psi(X|z)$ is the decoder, and $p(z)$ is standard Gaussian. The VAE is pretrained, and we fix its weights prior to training the other modules.

Point cloud policy. We train our point cloud policy with standard behavioral cloning (BC) loss, i.e. for the k -th skill, we sample a transition from the demonstration data, which contains the observed point clouds $\{P_i^o\}$, goal point clouds $\{P_j^g\}$, a tool point cloud P_k^{tool} , and the action of the tool a . Then, we match point clouds in the observation set to those in the goal set by finding the pairs of point clouds that are within a Chamfer Distance of ϵ : $\{(P_i^o, P_j^g) \mid D_{Chamfer}(P_i^o, P_j^g) < \epsilon\}$ and filter out the non-relevant point clouds. Last, we pass the filtered point clouds into the policy and minimize the following loss:

$$\mathcal{L}_{\pi_k} = \mathbb{E} [\|a - \pi_k(\{P_i^o\}, P_k^{tool}, \{P_j^g\})\|^2] \quad (3)$$

Feasibility predictor. We train the feasibility predictor for the k -th skill by regressing to the ground-truth feasibility label using mean squared error (MSE) as loss, i.e.

$$\mathcal{L}_{f_k} = \mathbb{E} \left[\left(f_k(\hat{U}^o, \hat{U}^g) - \mathbb{1}\{\hat{U}^o, \hat{U}^g \text{ is a positive pair}\} \right)^2 \right] \quad (4)$$

During training, we obtain positive pairs for the feasibility predictor by sampling two point clouds (P^{obs}, P^{goal}) from the same trajectory in the demonstration set. To find \hat{U}^o, \hat{U}^g , we first cluster the

Parameter	LiftSpread	GatherMove	CutRearrange	CRS + CRS-Twice
Yield stress	200	200	150	150
Ground friction	1.5	1.5	0.5	0.5
Young’s modulus (E)	5e3	5e3	5e3	5e3
Poisson’s ratio (ν)	0.15	0.15	0.15	0.15

Table 2: Parameters for simulation dough

138 observation and goal point clouds into two sets $\{P_i^o\}, \{P_j^g\}$ respectively. Then, we match point
 139 clouds in the observation set to those in the goal set by finding the pairs of point clouds that are within
 140 a Chamfer Distance of ϵ : $\{(P_i^o, P_j^g) \mid D_{Chamfer}(P_i^o, P_j^g) < \epsilon\}$. We then remove these point clouds
 141 from the corresponding set, since these are the point clouds that have already been moved to the target
 142 location in the goal. We can then encode the remaining point clouds into \hat{U}^o, \hat{U}^g using our VAE.

143 **Cost predictor.** We train the cost predictor by simply regressing to the Chamfer Distance (CD)
 144 between two entities represented by their latent vectors, i.e.

$$\mathcal{L}_c = \mathbb{E} \left[(c(\phi(P_i), \phi(P_j)) - D_{Chamfer}(P_i, P_j))^2 \right] \quad (5)$$

145 where P_i and P_j are point clouds of a single entity sampled from the dataset and ϕ is the encoder.
 146 There are two reasons that we train a cost predictor on latent vectors instead of directly computing
 147 the Chamfer Distance between two point clouds. For one, decoding each latent vector would greatly
 148 bottleneck the planning speed. Experiments on CutRearrange show that with our learned cost
 149 predictor, the planning takes 35s; on the other hand, if we decode the latent vectors and use the
 150 Chamfer Distance, even with a subsampled point cloud of 200 points, the planning takes 37200s
 151 (around 10 hours), which is impractical to use. Moreover, using a cost predictor can also offer us the
 152 flexibility to incorporate more complex reward functions in the future.

153 Finally, We train our policy, feasibility predictor, and cost predictor with the following loss:

$$\mathcal{L}_{PASTA} = \sum_{k=1}^K \mathbb{E} [\lambda_\pi \mathcal{L}_{\pi_k} + \lambda_f \mathcal{L}_{f_k} + \lambda_c \mathcal{L}_c] \quad (6)$$

154 We use $\lambda_\pi = 1$, $\lambda_f = 10$, and $\lambda_c = 1$ for all of our experiments.

155 2 Details on Simulation Experiments

156 2.1 Hyperparameters for Simulation Dough

157 We use PlasticineLab [5] for evaluating our simulation experiments. We provide the hyperparameters
 158 that are relevant to the properties of the dough in simulation to enhance the replicability of our results.
 159 See Table 2 for details.

160 2.2 Hyperparameters for DBSCAN

161 To cluster a point cloud, we use Scikit-learn’s [6] implementation of DBSCAN [7] with
 162 `eps=0.03`, `min_samples=6`, `min_points=10` for all of our environments. Further, we assign
 163 each noise point identified by DBSCAN to its closest cluster.

164 2.3 Hyperparameters for PASTA

165 Table 3 shows the hyperparameters used for PASTA in our simulation tasks. Planning for CRS-Twice
 166 requires a large amount of samples. Therefore, we modify the planner to improve sample efficiency.
 167 See Sec. 2.4 for details.

Training parameters	LiftSpread	GatherMove	CutRearrange	CRS	CRS-Twice
<i>Point Cloud VAE</i>					
learning rate	2e-3	2e-3	2e-3	2e-3	2e-3
latent dimension	2	2	2	2	2
<i>Feasibility predictor</i>					
learning rate	1e-4	1e-3	1e-4	1e-4	1e-4
batch size	256	256	256	256	256
noise on shape encoding σ_z	0	0	0	0.02	0.02
noise on position σ_t	0.01	0.01	0.005	0.01	0.01
<i>Cost predictor</i>					
learning rate	1e-4	1e-3	1e-4	1e-4	1e-4
batch size	256	256	256	256	256
<i>Policy</i>					
learning rate	1e-4	1e-3	1e-4	1e-4	1e-4
batch size	10	10	10	10	10
noise on point cloud	0.005	0.005	0.005	0.005	0.005
Planning parameters	LiftSpread	GatherMove	CutRearrange	CRS	CRS-Twice
learning rate	0.01	0.01	0.01	0.01	0.01
number of iterations	200	100	100	200	300
number of samples	5000	5000	5000	50000	500000

Table 3: Summary of hyperparameters used in PASTA. For CRS-Twice, we use the same model as CRS but modify the planner to have better sample efficiency.

168 2.4 Receding Horizon Planning for CRS-Twice

169 As the planning horizon increases, the number of possible skill sequences as well as the number of
170 possible attention structures increases exponentially. The task of CRS-Twice has a planning horizon
171 of 6 and is a much more difficult task to solve. As such, for this task, we specify the skill sequences
172 and use Receding Horizon Planning (RHP). Starting from the first time step, we follow the procedure
173 in Algorithm 1 but only optimize for H_{RHP} steps into the future and compare the achieved subgoal
174 with the final target to compute the planning loss. After optimization, we take the first subgoal from
175 the plan and discard the rest of the plan. We then repeat this process until we reach the overall
176 planning horizon H . In our experiments, we use $H_{RHP} = 3$. While we can perform model predictive
177 control and execute the first step before planning for the second step, we find this open-loop planning
178 and execution to be sufficient for the task.

179 3 Details on Real World Experiments

180 3.1 Heuristic Policies

181 Transferring the learned policies from simulation to the real world can be more difficult than transfer-
182 ring the planner itself, as the policies are affected more by the sim2real gap, such as the difference
183 in friction and properties of dough in the real world. To sidestep this challenge, for our real world
184 experiments, we design three heuristic policies: cut, push, and roll to execute the generated plans in
185 the real world.

186 Just like our learned policies in simulation, each heuristic policy takes in the current observation and
187 the generated subgoal in point clouds and outputs a sequence of desired end effector positions used
188 for impedance control. In addition, each policy takes in the attention mask provided by the planner
189 indicating the components of interest. The same DBSCAN procedure is used for this. The cut policy
190 first calculates the cutting point by computing the length ratio of the resulting components. Then it
191 cuts the dough and separates it such that the center of mass of each resulting component matches the
192 one in the subgoal. For the push policy, given a component and a goal component, the policy pushes
193 the dough in the direction that connects the two components' center of mass. The roll policy first

194 moves the roller down to make contact with the dough. Then, based on the goal component’s length,
 195 the policy calculates the distance it needs to move the roller back and forth when making contact
 196 with the dough.

197 3.2 Procedure for Resetting the Dough

198 To compare different methods with the same initial and target configurations, we first use a 3D-printed
 199 mold to fit the dough to the same initial shape. We then overlay the desired initial location on the
 200 image captured by the top-down camera and place the dough at the corresponding location in the
 201 workspace to ensure different methods start from the same initial location.

202 3.3 Procedure for the Human Baseline

203 Following the same procedure in section 3.2, we first reset the dough to the initial configuration.
 204 Then, we overlay the goal point cloud on the image captured by the top-down camera. The overlay
 205 image is shown on a screen and presented to the human in real-time when the human is completing
 206 the task.

207 3.4 Procedure for Making the Real Dough

Material	Quantity(g)	Baker’s percentage(%)
Flour	300	100
Water	180	60
Yeast	3	1

Table 4: All-purpose dough recipe

208 We follow the recipe shown in Table 4 to make the real dough. Following the tradition of baking, we
 209 use the baker’s percentage, so that each ingredient in a formula is expressed as a percentage of the
 210 flour weight, and the flour weight is always expressed as 100%. First, we take 300 grams of flour, 3
 211 grams of yeast, and 180 grams of water into a basin. Then, we mix the ingredients and knead the
 212 dough for a few minutes. Next, we use a food wrap to seal the dough in the basin and put them in the
 213 refrigerator to let the dough rest for 4-5 hours. Finally, we take out the dough from the refrigerator
 214 and reheat it with a microwave for 30-60 seconds to soften it.

215 4 Additional Experiments

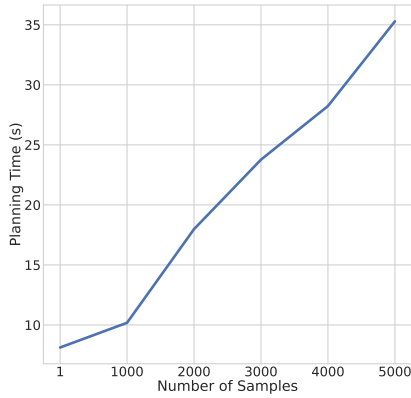
216 4.1 Ablation Studies

217 **Ablations on feasibility predictor.** Following
 218 the discussions in Sec 1.2, we train a feasibility
 219 predictor without adding any noise to show
 220 that adding noise helps with the optimization
 221 landscape during planning. We call this ablation
 222 *No Smoothing Feasibility*. As shown in Table 5,
 223 this variant only achieves half of the success rate
 224 of PASTA, suggesting the importance of noise
 225 during training.

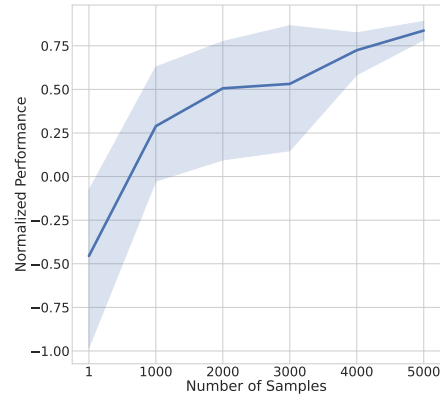
Ablation Method	Performance / Success
No Smoothing Feasibility	0.744 / 40%
Shared Encoder Policy	-0.304 / 0%
Tool Concat Policy	0.516 / 60%
Set without Filtering	0.360 / 20%
PASTA (Ours)	0.837 / 80%

Table 5: Additional ablation results from CutRearrange.

226 **Ablations on policy.** We consider two ablation
 227 methods for our set policy. First, we consider a *Shared Encoder Policy* that takes in the latent vectors
 228 from the encoder and uses a max pooling layer followed by an MLP to produce the action. The
 229 architecture is very similar to our Set Feasibility Predictor. Our results in Table 5 show that this
 230 architecture has zero success in our task. We hypothesize that this is because the entity encoding can



(a) Planning time v.s. number of samples



(b) Planning performance v.s. number of samples

Figure 2: Planning time and performance with varying number of samples in CutRearrange. We show the mean and standard deviation of performance over 5 runs.

231 be unstable during the skill execution. For example, during cutting, the dough slowly transitions from
 232 one piece to two pieces, making the input to the policy unstable.

233 Second, we compare with a *Tool Concat Policy* that takes in the observation and goal point cloud of
 234 the dough, passes them through a PointNet++ [3] encoder to produce a feature, and then concatenates
 235 the tool state to the feature. The concatenated feature is passed through a final MLP to output the
 236 action. In comparison, the set policy in PASTA takes the point cloud of the tool and concatenates
 237 it with the dough in the point cloud space before passing it to the PointNet. We hypothesize that
 238 this way allows PointNet to reason more easily about the spatial relationships between the tool point
 239 cloud and the dough point cloud. Results in Table 5 highlight the advantage of using a point cloud to
 240 represent the tool.

241 **Ablation on set representation.** We consider a variant of PASTA *Set without Filtering*, which uses
 242 the same set representation as PASTA, but does not filter entities that are approximately the same
 243 both during training and testing. This filtering is only possible with a set representation and we want
 244 to show the advantage of this filtering. For this ablation, during training, the feasibility predictor
 245 takes in all the entities in the scene in set representation, and the policy takes in the concatenation
 246 of point clouds from each entity. During planning, we do not enumerate attention structures but
 247 instead optimize for all the entities. As shown in Table 5, without filtering, this ablation performs
 248 significantly worse than PASTA, showing that filtering is an important advantage enabled by our set
 249 representation.

250 4.2 Visualization of the Latent Space

251 We visualize the latent space of PASTA in CutRearrange in Figure 3 and visualize the latent space
 252 of Flat 3D baseline in Figure 4 for comparison. Since we use a latent dimension of 2 for all of
 253 our environments, we can visualize the original latent space without applying any dimensionality
 254 reduction techniques. PASTA only encodes the shape of each entity and thus can better model the
 255 variations in shapes. On the other hand, Flat 3D couples the shape variation with the relative position
 256 of two entities. This makes a flat representation difficult to generalize compositionally to scenes with
 257 different numbers of entities or scenes with entities that have novel relative spatial locations to each
 258 other.

259 **4.3 Runtime of PASTA**

260 We implement the planning in the latent set representation in an efficient way, which can plan with
 261 multiple different structures in parallel on a GPU. To demonstrate the efficiency of PASTA, we vary the
 262 number of samples used for planning and record the planning time and final performance. We conduct
 263 the experiments in CutRearrange. Figure 2a shows that the planning time scales approximately
 264 linearly with the number of samples, and Figure 2b shows the planning performance versus the
 265 number of samples. As the result suggests, PASTA can achieve its optimal performance with a very
 266 short amount of planning time (under 1 minute) for the majority of our tasks. Finally, we summarize
 the planning time for all of our tasks in simulation in Table 6.

	LiftSpread	GatherMove	CutRearrange	CRS	CRS-Twice*
Planning time (seconds)	58	35	35	307	7810

Table 6: Summary of planning time of PASTA in all of the simulation tasks. CRS-Twice uses Receding Horizon Planning, which results in an increase in planning time.

267

268 **4.4 Additional metrics for real world experiments**

269 We also quantitatively computed the action error v.s. subgoal error for our real world trajectories. The
 270 results are shown in Table 7. From the results in the table, our planned goal is closer to the ground
 271 truth goal than the achieved goal, measured by the Earth Mover’s Distance (EMD), which shows that
 272 the controller does not compensate for the error of the planner.

	CutRearrange	CRS	CRS-Twice
EMD(planned goal, ground-truth goal)	0.038 ± 0.004	0.027 ± 0.004	0.029 ± 0.002
EMD(reached goal, ground-truth goal)	0.056 ± 0.007	0.044 ± 0.006	0.054 ± 0.016

Table 7: Action error v.s. subgoal error for real world experiments. For each task, the mean \pm std for 4 trajectories are shown.

273 **4.5 Robustness of PASTA**

274 We show that PASTA is robust to two types of variations and can retain high performance.

275 **Robust to planning horizon** First, we increase the planning horizon from the minimal length for
 276 the task (3) to twice the minimal length (6), and we observe that PASTA retains a high performance
 277 across all horizons. The results are shown in Table 8. This suggests that in practice, one can specify a
 278 maximum planning horizon for PASTA when the exact horizon is unknown.

279 **Robust to distractors** Second, we show that PASTA is robust to distractors in the scene. We add 2
 280 distractor objects in CRS (which makes the scene have 4 objects in total). We observe that PASTA
 281 retains a normalized performance of 0.879 and 100% success rate (without distractor: 0.896/100%)
 282 using the same amount of samples to plan. Our planner is able to ignore the distractors using our
 283 attention structure at every step to only attend to the relevant components in the scene. We also added
 284 an example trajectory with distractor dough pieces to our website under “CRS with distractors”.

Planning Horizon	3	4	5	6
Performance	0.896	0.866	0.90	0.878
Success Rate	5/5	4/5	5/5	4/5

Table 8: PASTA’s performance v.s. varying numbers of planning horizon in CRS.

285 5 Further Discussion on Limitations and Future Work

286 **More Efficient Planning** Planning skill sequences with a large search space is a challenging problem
287 by itself but much progress has been made by the task and motion planning community to obtain
288 a plan skeleton [8, 9, 10]. For example, Caelan et al. [8] propose two methods, the first one is to
289 interleave searching the skill sequence with lower-level optimization and the second one is to have
290 lazy placeholders for some skills. Danny et al. [10] propose to predict skill sequences from visual
291 observation. Recent works have also explored finding skill sequences using pre-trained language
292 models [11, 12].

293 **Sim2Real Transfer for Real Dough** One possible approach is to train with domain randomization
294 to make the policy more robust to changing dynamics (e.g. stickiness) of dough. Another option
295 is to perform online system identification of the dough dynamics parameters [13, 14] or real2sim
296 methods [15, 16]. In future work, we can also integrate our method with other works that perform
297 low-level dough manipulation in the real world, such as recent work from Qi et al. [17].

298 **Goal Specification** Our planner requires specifying the goal with a point cloud and coming up
299 with a point cloud goal is not always easy. However, rapid progress are being made with language-
300 conditioned manipulation and future work can combine language to specify more diverse tasks.

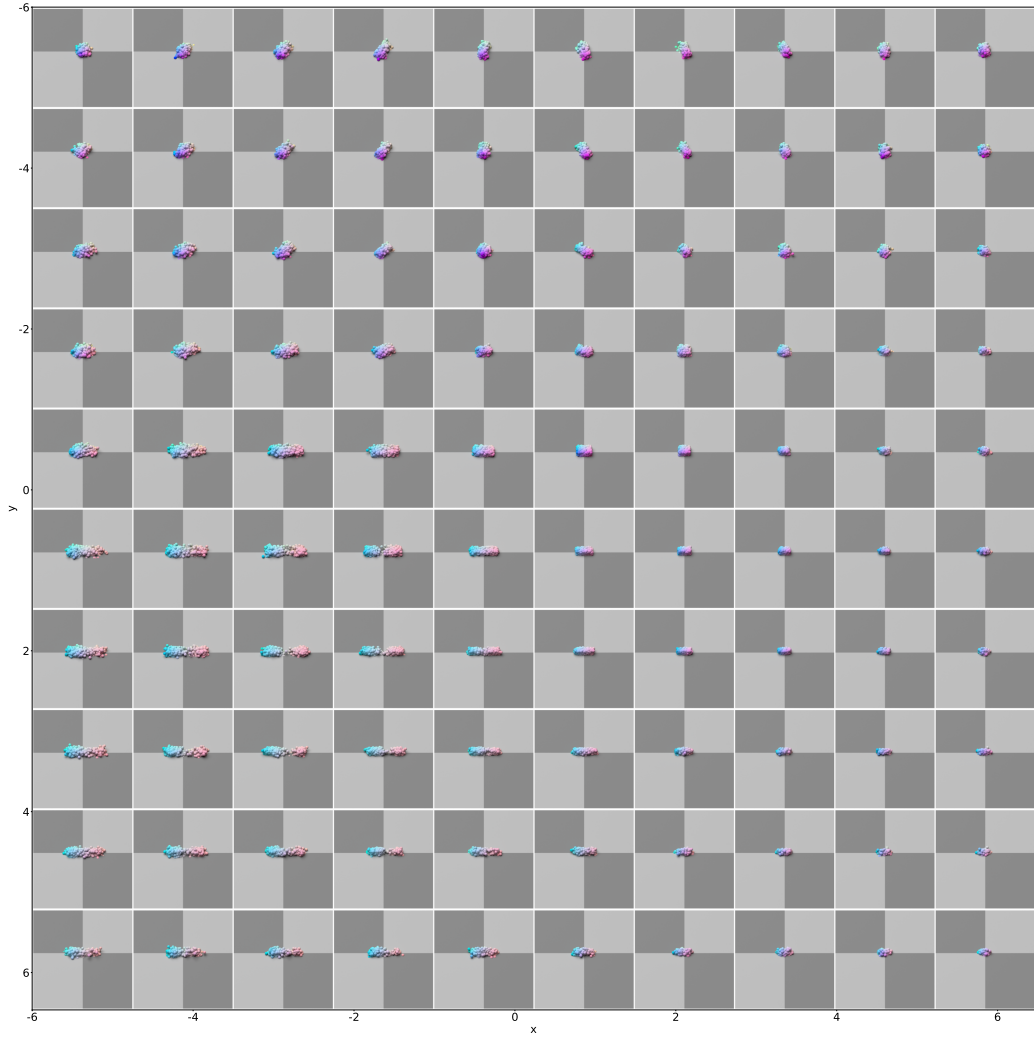


Figure 3: Latent space of PASTA in CutRearrange. We sample coordinates on a grid from the 2D latent space encoding the shapes and then decoding each latent vector into a point cloud. We then rearrange the decoded point cloud into the grid based on the corresponding coordinates in the latent space.

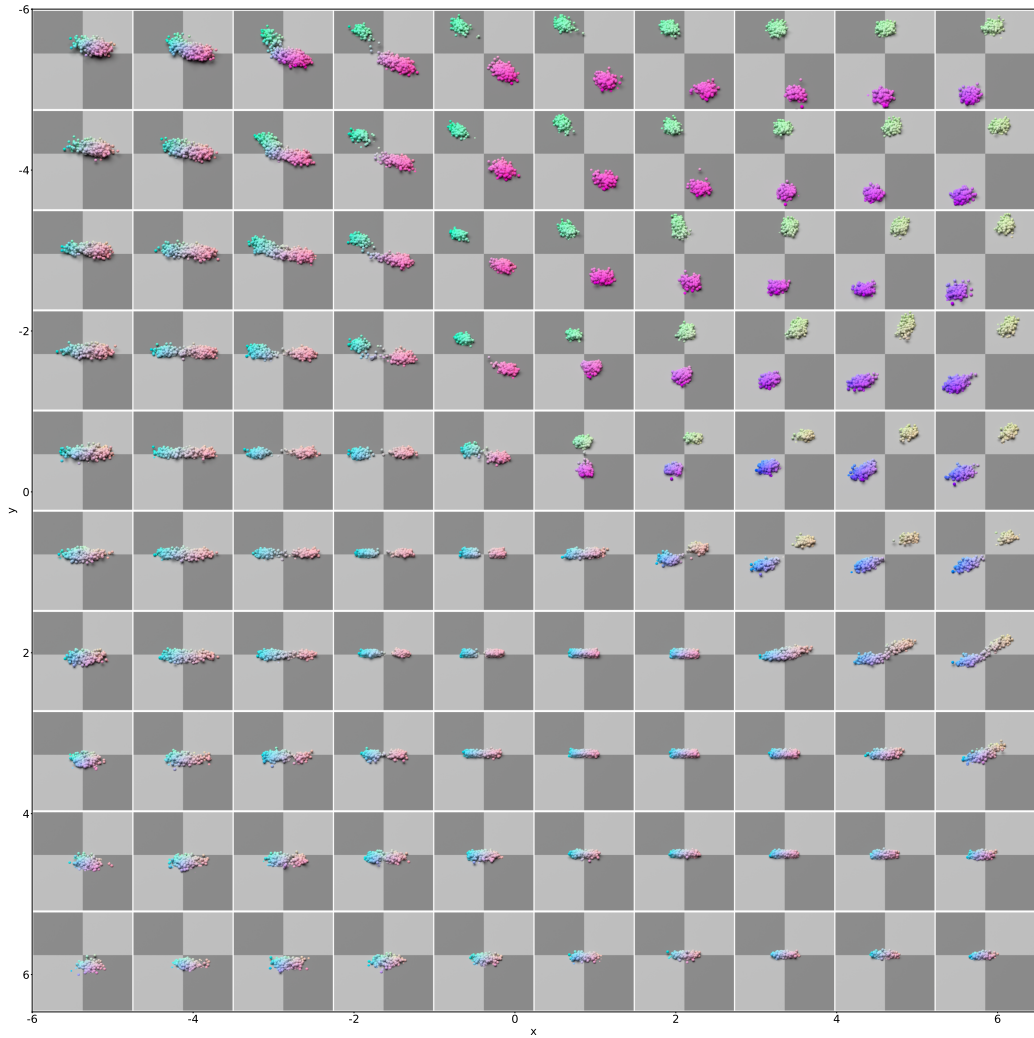


Figure 4: Latent space of Flat 3D in CutRearrange. We sample coordinates on a grid from the 2D latent space encoding the shapes and then decoding each latent vector into a point cloud. We then rearrange the decoded point cloud into the grid based on the corresponding coordinates in the latent space.

References

- [1] G. Yang, X. Huang, Z. Hao, M.-Y. Liu, S. Belongie, and B. Hariharan. Pointflow: 3d point cloud generation with continuous normalizing flows. *ICCV*, 2019.
- [2] X. Lin, Z. Huang, Y. Li, J. B. Tenenbaum, D. Held, and C. Gan. Diffskill: Skill abstraction from differentiable physics for deformable object manipulations with tools. *ICLR*, 2022.
- [3] C. R. Qi, L. Yi, H. Su, and L. J. Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. *Advances in neural information processing systems*, 30, 2017.
- [4] M. Fey and J. E. Lenssen. Fast graph representation learning with PyTorch Geometric. In *ICLR Workshop on Representation Learning on Graphs and Manifolds*, 2019.
- [5] Z. Huang, Y. Hu, T. Du, S. Zhou, H. Su, J. B. Tenenbaum, and C. Gan. Plasticinelab: A soft-body manipulation benchmark with differentiable physics. In *International Conference on Learning Representations*, 2020.
- [6] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- [7] M. Ester, H.-P. Kriegel, J. Sander, X. Xu, et al. A density-based algorithm for discovering clusters in large spatial databases with noise. In *kdd*, volume 96, pages 226–231, 1996.
- [8] C. Garrett, T. Lozano-Pérez, and L. Kaelbling. Sample-based methods for factored task and motion planning. 2017.
- [9] B. Kim and L. Shimanuki. Learning value functions with relational state representations for guiding task-and-motion planning. In *Conference on Robot Learning*, pages 955–968. PMLR, 2020.
- [10] D. Driess, J.-S. Ha, and M. Toussaint. Deep visual reasoning: Learning to predict action sequences for task and motion planning from an initial scene image. *arXiv preprint arXiv:2006.05398*, 2020.
- [11] M. Ahn, A. Brohan, N. Brown, Y. Chebotar, O. Cortes, B. David, C. Finn, K. Gopalakrishnan, K. Hausman, A. Herzog, et al. Do as i can, not as i say: Grounding language in robotic affordances. *arXiv preprint arXiv:2204.01691*, 2022.
- [12] W. Huang, P. Abbeel, D. Pathak, and I. Mordatch. Language models as zero-shot planners: Extracting actionable knowledge for embodied agents. *arXiv preprint arXiv:2201.07207*, 2022.
- [13] W. Yu, J. Tan, C. K. Liu, and G. Turk. Preparing for the unknown: Learning a universal policy with online system identification. *arXiv preprint arXiv:1702.02453*, 2017.
- [14] A. Kumar, Z. Fu, D. Pathak, and J. Malik. Rma: Rapid motor adaptation for legged robots. *arXiv preprint arXiv:2107.04034*, 2021.
- [15] F. Ramos, R. C. Possas, and D. Fox. Bayessim: adaptive domain randomization via probabilistic inference for robotics simulators. *arXiv preprint arXiv:1906.01728*, 2019.
- [16] Y. Chebotar, A. Handa, V. Makoviychuk, M. Macklin, J. Issac, N. Ratliff, and D. Fox. Closing the sim-to-real loop: Adapting simulation randomization with real world experience. In *2019 International Conference on Robotics and Automation (ICRA)*, pages 8973–8979. IEEE, 2019.
- [17] C. Qi, X. Lin, and D. Held. Learning closed-loop dough manipulation using a differentiable reset module. *IEEE Robotics and Automation Letters*, 7(4):9857–9864, 2022.