

1 Appendix

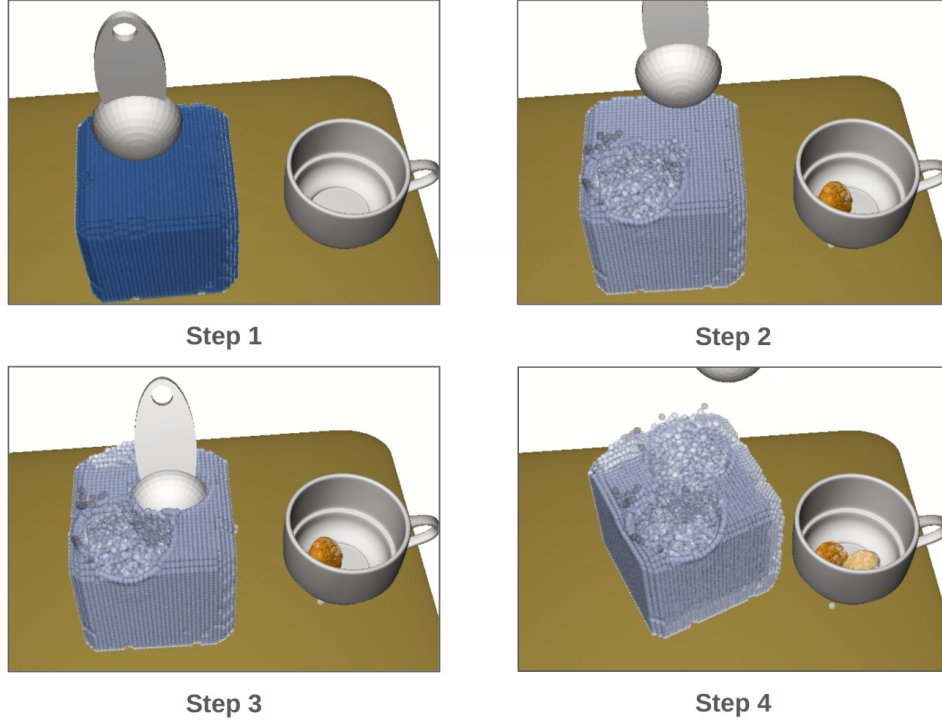


Figure 1: Ice-Cream Scooping Task. (Step 1) Initial setup of the scene, with the spoon positioned above the ice cream cube. A white cup is placed beside the ice cream block to receive the scooped pieces. (Step 2) After the first scoop, a distinct ice cream piece appears in the cup, shown in a different color to indicate it is a separate connectivity cluster. (Step 3) The spoon performs a second scooping, inserting into the ice cream cube to gather another piece. (Step 4) The task is completed with two distinct ice cream pieces, each in its own cluster, now present in the cup.

2 1.1 Additional Manipulation Tasks

While our main focus is on multi-step cutting, our simulator, topology discovery mechanism, and spectral reward function extend naturally to other deformable-object tasks. Below, we describe two illustrative examples that highlight the potential for extending our framework, using enumerated steps for each task’s execution.

7 1.1.1 Cream Writing

In this task, the objective is to “write” a target word (e.g., CORL) by extruding a soft, von Mises-plastic “cream” material onto a flat surface. The cream is modeled in MPM with low yield stress and high plasticity, allowing it to retain its shape after extrusion. The agent controls a “cream pen” at a fixed height above the table. Its state is represented by its (x, y) position, and the action is a 3-dimensional vector $(\Delta x, \Delta y, b)$, where $(\Delta x, \Delta y)$ specifies the pen movement and $b \in \{0, 1\}$ toggles the extrusion.

During the writing process, each continuous stroke formed while $b = 1$ is considered a single connectivity cluster. If the pen is lifted and re-positioned, the next stroke is assigned to a new cluster. Thus, the clustering mechanism is naturally defined by the writing process, not by explicit color control. The visualization in Figure 2 uses different colors to indicate separate clusters, a distinction that becomes critical when computing the spectral-based reward.

Execution steps:

1. Position the pen at the starting point for the letter “C,” toggle $b = 1$, and move rightward to trace the stroke, forming the first connectivity cluster.

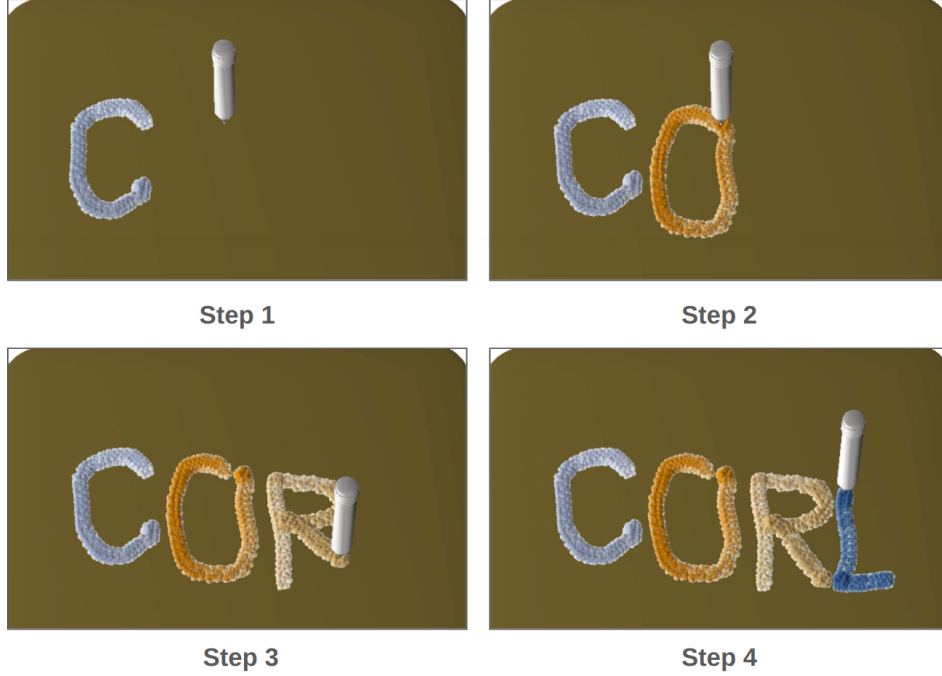


Figure 2: Cream Writing Task. (Step 1) The cream pen extrudes a continuous stroke to form the letter “C.” (Step 2) The pen lifts, repositions, and extrudes the letter “O.” (Step 3) Next, it draws “R.” (Step 4) Finally, it writes “L.” Each letter is rendered in a different color to indicate that they belong to separate connectivity clusters. The color differentiation is not manually controlled but automatically generated based on the clustering, which is leveraged during the spectral-based reward computation.

- 22 2. Lift the pen, move to the position for the letter “O,” toggle $b = 1$, and extrude to form the second
- 23 cluster.
- 24 3. Reposition to the third slot, toggle $b = 1$, and draw the letter “R,” creating a new cluster.
- 25 4. Finally, move to the fourth slot, toggle $b = 1$, and write the letter “L,” forming the final cluster.

26 Our spectral reward function evaluates how well each cluster corresponds to the intended letter shape.

27 This clustering mechanism provides a natural way to segment distinct components, allowing the

28 reward function to assess both spatial structure and connectivity. However, the sequential planning

29 of strokes remains a complex decision-making challenge, distinct from the cutting task, and thus

30 extending PDDP to handle such “additive” tasks is a promising direction for future work.

31 1.1.2 Ice-Cream Scooping

32 In this task, the objective is to scoop a chunk from a block of “pudding”-like ice cream (elastoplastic

33 MPM) and deposit it into a cup. The agent controls a spoon at a fixed height, using a 2D (x, y) action

34 space to position the scoop. Once positioned, the agent executes a predefined scooping primitive as

35 follows:

36 Execution steps:

- 37 1. Lower the spoon into the ice cream block to initiate contact.
- 38 2. Translate the spoon forward to penetrate and gather the material.
- 39 3. Lift the spoon while rotating it to maintain a horizontal orientation, preventing spillage.
- 40 4. Move the spoon above the cup.
- 41 5. Tilt or open the spoon to release the scooped chunk into the cup.

42 The scooped chunk forms a distinct connectivity cluster, separate from the remaining block. As

43 illustrated in Figure 1, each scooped piece is automatically assigned a different color to indicate it

44 belongs to a separate cluster. This distinction is crucial when computing our spectral-based reward,

45 as it allows for quantifying successful segmentation and accurate placement of the scooped pieces.

1.1.3 Discussion

These two examples demonstrate the extensibility of our MPM-based environment, topology discovery, and spectral reward function to a broader range of deformable manipulation tasks beyond multi-step cutting. The use of distinct clusters to represent separate connectivity components is a key mechanism for computing spectral-based rewards, applicable in both “additive” and “scooping” tasks. Designing specialized action primitives and learning-based planners for these scenarios remains a promising direction for future exploration.

1.2 Pyramid Cutting Task

Motivation and Objective. – To further evaluate the generalization capability of our spectral reward function to more complex and arbitrary goal shapes, we introduce the pyramid cutting task (Figure 3). In contrast to the predefined slicing, sticking, and dicing tasks, this experiment requires the agent to carve out a corner segment from a cubic block to produce a pyramid-shaped fragment. The target shape is defined as a triangular pyramid with sloped surfaces, representing a more intricate and asymmetrical geometry than previous tasks.

The objective of this experiment is to assess whether the MPPI planning method, driven solely by our spectral reward function, can effectively discover a plausible cutting trajectory to achieve the desired pyramid shape without any task-specific tuning or retraining. This setup challenges the reward function to guide cutting actions that align with more complex, multi-faceted surfaces, testing its robustness and generalizability.

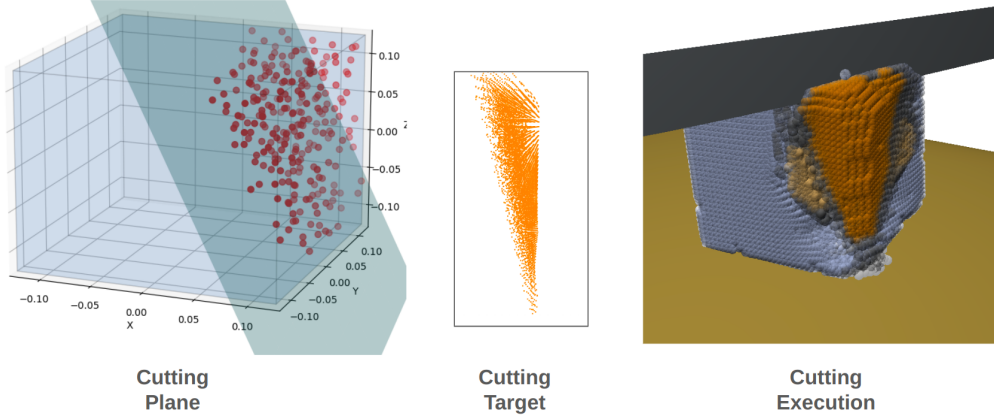


Figure 3: Pyramid Cutting Task. (Left) Cutting Plane: Visualization of the cutting plane planned by the MPPI controller based on our spectral reward function. The plane is strategically positioned to carve out a pyramid-shaped fragment from the corner of the cube. (Middle) Cutting Target: The target goal shape, shown as a point cloud representation, highlighting the desired pyramid structure that serves as a guide for the cutting task. (Right) Cutting Execution: The scene depicts the execution of the planned cutting trajectory, where the knife removes the corner section to achieve the specified pyramid shape. The orange particles indicate the material identified as the target segment to be removed, aligning with the planned cutting plane.

Takeaway. – The results of the pyramid cutting task clearly demonstrate that our spectral reward function effectively generalizes to more intricate geometric goals, even in cases involving angled and asymmetric surfaces. As shown in Figure 3, the MPPI planner successfully identifies a cutting plane that produces a distinct pyramid-shaped fragment from the cube’s corner, aligning well with the target structure. The qualitative alignment between the planned cutting plane and the final extracted segment validates that the reward function remains a reliable guidance signal, even for complex and non-standard shapes. This outcome highlights the adaptability of our framework, suggesting its potential to handle more sophisticated and arbitrary cutting tasks in future applications.

73 1.3 Perception Model Architecture Details

74 In practice, we design two distinct styles of perception models to explore their effectiveness in
75 our setup: a graph-based model and a point-based model. The graph-based model leverages graph
76 convolutional networks (GCNs) and joint attention mechanisms to process structured graph data,
77 while the point-based model employs a hierarchical PointNet-style architecture using SPALPA blocks
78 to handle point cloud data. After extensive evaluation, we adopt the graph-based model as our final
79 choice. This decision is driven by the significant reduction in parameter count (approximately four
80 times fewer parameters) and faster training and inference speed compared to the point-based model.
81 The PointNet-based architecture, though capable of capturing finer local features, incurs substantial
82 computational overhead due to its deeper SPALPA structure and multi-scale processing layers.

83 1.3.1 PointNet-Based Perception Model Architecture

84 The PointNet-based architecture processes point cloud data through hierarchical SPALPA blocks for
85 multi-scale feature extraction. The key components are:

- 86 • **Topological Encoder:** The primary perception encoder in this architecture, responsible for ex-
87 tracting features from 32-dimensional topological inputs. It comprises 5 SPALPA blocks that
88 progressively increase feature dimensions as follows:

$$96 \rightarrow 192 \rightarrow 384 \rightarrow 768 \rightarrow 1536$$

89 Each SPALPA block includes:

- 90 – **Local Attention:** Extracts localized spatial relationships through ‘Conv2d’.
 - 91 – **Global Attention:** Aggregates broader context using multi-scale convolutional layers.
 - 92 – **Grouper:** Implements ‘QueryAndGroup’ for spatial neighborhood aggregation.
 - 93 • **Action Encoder:** Encodes 2-dimensional action inputs using the same hierarchical structure as the
94 Topological Encoder, ensuring feature alignment and consistency.
 - 95 • **Cross-Attention Network:** Integrates topological and action features through CrossAttnNet
96 modules, employing PreNorm Attention, cross-attention, and gated residual connections.
 - 97 • **Decoder:** Reconstructs the feature maps and propagates them to the original point cloud resolution
98 using ‘FeaturePropagation’ blocks, merging multi-scale features.
 - 99 • **Output Heads:**
 - 100 – **Segmentation Head:** Projects features to segmentation classes using ‘Conv1d’ layers.
 - 101 – **Point Cloud Head:** Predicts point coordinates through residual blocks and a final ‘Linear’ layer.
- 102 Despite its comprehensive feature extraction capability, the PointNet-based model incurs con-
103 siderable computational overhead, resulting in slower training and inference compared to the
104 graph-based model.

105 1.3.2 Graph Network Architecture

106 The Graph Network architecture processes structured graph data using GCN layers and attention-
107 based mechanisms. The core perception module in this architecture is the **t_graph_encoder**, which
108 encodes topological information through GCN layers.

- 109 • **t_graph_encoder:** The designated perception module in this architecture, processing node and
110 label embeddings through GCN layers structured as:

$$3 \rightarrow 96 \rightarrow 96$$

111 Key submodules include:

- 112 – **Node Encoder:** Projects 3-dimensional node features to a 96-dimensional latent space.
- 113 – **Label Encoder:** Processes 32-dimensional node labels through GCN layers.
- 114 – **Embedding Encoder:** Integrates 192-dimensional precomputed embeddings into a unified
115 96-dimensional space.
- 116 • **a_graph_encoder:** Processes action-related graph data using a similar structure to the
117 ‘t_graph_encoder’, but with 2-dimensional action inputs.
- 118 • **Joint Graph Transformer:** Integrates features from both graph encoders through multi-layer
119 attention and cross-attention modules:

$$96 \rightarrow 64 \rightarrow 32 \rightarrow 16 \rightarrow 8 \rightarrow 4$$

Table 1: Empirical Studies of Pretraining Strategies for Perception Models: We evaluate the effectiveness of graph-based and point-based perception models under various pretraining strategies, including sorting the input point cloud, applying Gaussian noise for data augmentation, and downsampling the input points before processing. The final model selected for deployment is the Graph-Based model with sorting, Gaussian noise, and 256 downsampled points, balancing accuracy and computational efficiency.

Perception Model	With Sort	Gaussian Noise	Downsample Points	Seen Accuracy	Unseen Accuracy
Graph-Based	✗	✗	128	66.4	48.7
Graph-Based	✓	✗	128	73.8	52.1
Graph-Based	✓	✓	128	75.4	69.3
Graph-Based (Final)	✓	✓	256	77.1	74.1
PointNet-Based	✗	✓	2048	73.2	61.2
PointNet-Based	✓	✓	512	79.2	75.4

- **Query Point Encoder:** Processes query nodes in the graph structure through ‘Linear’ and ‘ReLU’ layers:

$$3 \rightarrow 64 \rightarrow 96$$

• Output Heads:

- **Node Position Head:** Outputs node coordinates through residual blocks.
- **Node Feature Head:** Projects node features to a 32-dimensional space.

The graph-based model provides a more parameter-efficient and computationally feasible architecture while maintaining robust feature extraction and representation capabilities, making it the preferred choice in our implementation.

Summary: – We explore both PointNet-based and Graph Network architectures to identify the optimal perception model for our setup. While the PointNet-based model employs extensive SPALPA blocks for multi-scale feature extraction, its computational cost and parameter count are substantially higher than the graph-based model. Consequently, we adopt the graph-based architecture, leveraging the **t_graph_encoder** as the primary perception module due to its significantly lower parameter count, faster training and inference speeds, and effective integration of topological and action features.

1.3.3 Empirical Evaluation of Perception Models

Discussion and Model Selection: – Table 1 presents a comparative evaluation of graph-based and point-based perception models under various pretraining strategies. The primary factors analyzed include sorting based on the x-axis, Gaussian noise for data augmentation, and the degree of point cloud downsampling.

- **Sorting and Gaussian Noise:** Incorporating sorting and Gaussian noise consistently improves model accuracy for both seen and unseen data. For the graph-based model, enabling both sorting and Gaussian noise increased unseen accuracy from 48.7% to 69.3% under a 128-point downsampling setting.
- **Downsampling Strategy:** The impact of point cloud resolution is evident as increasing the number of points from 128 to 256 in the graph-based model further improved unseen accuracy to 74.1%. This result underscores the importance of maintaining sufficient point density to preserve critical spatial information.
- **Comparison of Models:** The point-based model with 512 points achieves 75.4% unseen accuracy, slightly outperforming the graph-based model with 256 points. However, the computational overhead associated with the point-based model is significantly higher, consuming approximately four times more parameters and resulting in slower training and inference.

Final Model Selection: – Based on the empirical findings, the graph-based model with sorting, Gaussian noise, and 256 downsampled points is selected as the final perception model for deployment. This configuration achieves a balanced trade-off between computational efficiency and segmentation accuracy, making it a practical choice for large-scale robotic manipulation tasks.

1.4 Policy Model Architecture Details

In our framework, the policy model architecture is structured to handle skill-conditioned action generation through a structured diffusion process. The model, named **PDDP (Particle-Based Diffusion Policy)**, leverages adaptive normalization, timestep embeddings, and modular diffusion blocks to effectively model action trajectories. Notably, the observation encoder in PDDP utilizes the perception encoder previously defined in the **Perception Model Architecture** section, ensuring consistent feature extraction and representation across the pipeline.

1.4.1 Policy Model Architecture - PDDP

The PDDP architecture is structured as follows:

1. Encoders: –

- **Observation Encoder:** The observation encoder adopts the perception encoder design outlined previously, specifically utilizing the graph-based **t_graph_encoder** as the primary perception module. This encoder processes topological node features through a GCN and includes three key submodules:

- **Node Encoder:** Projects 3-dimensional node features to a 96-dimensional latent space through a 2-layer GCN:

$$3 \rightarrow 96 \rightarrow 96$$

- **Label Encoder:** Transforms 32-dimensional node labels to a 96-dimensional space using the same GCN structure.

- **Embedding Encoder:** Integrates 192-dimensional precomputed embeddings, projecting them to 96 dimensions.

The outputs from each encoder component are processed through a shared ‘LayerNorm’ to stabilize feature representation.

- **Goal Encoder:** Encodes 3-dimensional goal vectors through a fully connected network:

$$3 \rightarrow 64 \rightarrow 96$$

ReLU activation is applied to the intermediate layer.

- **Timestep Embedding (sigma_map):** Encodes the diffusion timestep using a 2-layer MLP:

$$256 \rightarrow 128 \rightarrow 128$$

A SiLU activation is applied after the first linear layer to maintain smooth gradient flow.

- **Action History Layer:** Integrates action history features through a fully connected layer:

$$21 \rightarrow 128$$

- **Conditional Layer:** Aggregates timestep, goal, and action history embeddings into a unified feature space:

$$256 \rightarrow 128$$

2. Backbone - Diffusion Blocks (DDiTBlock): –

The core processing backbone of PDDP consists of a sequence of 12 **DDiTBlocks**. Each DDiTBlock contains:

- **Self-Attention Layer:** Applies self-attention to the feature space with 96-dimensional query, key, and value projections:

$$96 \rightarrow 288 \rightarrow 96$$

This allows the model to capture dependencies across all points in the point cloud, facilitating information exchange across nodes.

- **MLP Layer:** Implements a 2-layer MLP with GELU activation:

$$96 \rightarrow 384 \rightarrow 96$$

This structure refines the feature representation after the attention operation, maintaining non-linear feature transformations.

- **Dropout:** Applied to both the attention and MLP layers with a probability of 0.4 to mitigate overfitting and stabilize training.
- **Adaptive LayerNorm (adaLN):** Adaptive LayerNorm is applied to each block using modulation inputs derived from both the timestep embedding and goal encoder. The modulation layer is structured as:

$$128 \rightarrow 576$$

This mechanism enables each DDiTBlock to dynamically adjust its feature scaling based on task-specific conditions.

3. Output Layer: –

The final output layer processes the feature representation generated by the diffusion blocks to predict binary classification logits for each point in the point cloud. The output shape is defined as:

$$\text{Output Shape: } [\text{batch size}, \text{num points}, 2]$$

A linear projection is applied to transform the 96-dimensional feature space into 2-dimensional logits, representing the binary classification output for each point. This design aligns the output structure with the expected action space while maintaining consistency across the diffusion layers.

Summary: – The PDDP architecture is structured to leverage a modular pipeline consisting of encoders for observation, goal, and timestep embeddings, a backbone of 12 DDiTBlocks with adaptive normalization and self-attention mechanisms, and a binary classification output layer for per-point action prediction. The use of the **t_graph_encoder** as the primary perception module ensures consistent feature extraction and integration across both the perception and policy networks, promoting robust action generation in diverse manipulation tasks.

1.5 Material Point Method in Details

The Material Point Method (MPM) is a hybrid Lagrangian–Eulerian scheme originally introduced by Sulsky *et al.* for solid mechanics. In MPM, state variables (mass, momentum, deformation) are carried on material points (particles), while a fixed Eulerian grid is used to solve the equations of motion. This split representation combines the mesh-free advantages of particle methods with the stability and boundary-handling capabilities of grid methods.

Governing Equations – MPM begins from the continuum balance laws in Eulerian form:

$$\frac{D\rho}{Dt} + \rho \nabla \cdot v = 0, \quad \rho \frac{Dv}{Dt} = \nabla \cdot \sigma + \rho g,$$

where ρ is the mass density, v the velocity field, σ the Cauchy stress tensor, and g the body-force (e.g. gravity). The material derivative D/Dt captures convection of field quantities with the flow.

Weak Form and Discretization – To derive a tractable discretization, one multiplies the momentum equation by a test function $q(x)$ and integrates over the domain Ω^n at time step n . Integration by parts moves spatial derivatives onto q , yielding the weak form. We then partition Ω^n into particle subdomains Ω_p^n , and approximate both trial and test functions using B-spline shape functions $N_i(x)$ centered at grid nodes i . This Galerkin-style projection leads to discrete nodal equations that are assembled via sums over particles.

Particle-to-Grid (P2G) Transfer – Each particle p holds:

$$\{m_p, V_p, v_p, F_p\}$$

for mass m_p , volume V_p , velocity v_p , and deformation gradient F_p . To project onto the grid, we compute at each node i :

$$m_i = \sum_p m_p N_i(x_p), \quad (mv)_i = \sum_p m_p v_p N_i(x_p), \quad f_i^{\text{int}} = - \sum_p V_p \sigma_p \nabla N_i(x_p),$$

where σ_p is obtained from the chosen constitutive law (e.g. hyperelastic) evaluated at F_p . These transfers ensure exact conservation of mass and momentum.

Grid Update – On the Eulerian grid, we update nodal momentum via a symplectic (explicit) Euler step:

$$(mv)_i^{n+1} = (mv)_i^n + \Delta t(f_i^{\text{int}} + m_i g), \quad v_i^{n+1} = \frac{(mv)_i^{n+1}}{m_i}.$$

This step advances velocities under both internal stresses and external forces, while preserving stability for moderate time steps.

Grid-to-Particle (G2P) Transfer – After updating the grid, we interpolate back to particles:

$$v_p^{n+1} = \sum_i N_i(x_p) v_i^{n+1}, \quad C_p = \sum_i v_i^{n+1} \nabla N_i(x_p)^T,$$

where C_p is an affine velocity gradient matrix (used in variants like APIC) that captures sub-cell velocity variation.

Deformation Gradient Update – The particle deformation gradient evolves according to

$$\frac{dF}{dt} = (\nabla v) F.$$

Using the interpolated velocity gradient C_p , we discretize:

$$F_p^{n+1} = (I + \Delta t C_p) F_p^n,$$

so that F_p accumulates the local deformation history on each particle.

Stress Computation – For hyperelastic materials, one defines a strain-energy density $\Psi(F)$ and computes the first Piola–Kirchhoff stress

$$P_p = \frac{\partial \Psi}{\partial F}(F_p),$$

then the corresponding Cauchy stress

$$\sigma_p = \frac{1}{\det F_p} P_p F_p^T.$$

This stress is used in the P2G transfer to produce internal forces.

Algorithm Summary – At each timestep, MPM executes:

1. **P2G:** Transfer $\{m_p, v_p, F_p\}$ to grid nodes $\{m_i, v_i, f_i\}$.
2. **Grid Update:** Integrate nodal momentum, compute v_i^{n+1} .
3. **G2P:** Interpolate v_i^{n+1} and ∇v back to particles.
4. **State Update:** Update each F_p^{n+1} and compute σ_p .
5. **Advection:** Move particles: $x_p^{n+1} = x_p^n + \Delta t v_p^{n+1}$.
6. **Reset:** Clear grid variables for the next iteration.

Thanks to its hybrid nature, MPM can robustly simulate extreme deformations, fracture propagation, and multi-body contact without remeshing, making it a powerful tool in graphics, engineering, and robotic manipulation contexts.

1.6 Particle-based Damage-tracking and Topological Reconstruction in Details

We propose a robust particle-based method to accurately determine whether a cutting action successfully separates an object into discrete pieces. Our approach integrates particle-level damage tracking with topological surface reconstruction using the Material Point Method (MPM) as shown in Figure 4. Within our MPM framework, each particle’s deformation state is computed, and we define a particle as “damaged” based on critical compression and stretch thresholds in its deformation gradient \mathbf{F} . Formally, a particle p is classified as damaged if:

$$J_p \leq (1 - \epsilon_c)^m \quad \text{or} \quad J_p \geq (1 + \epsilon_s)^m, \quad (1)$$

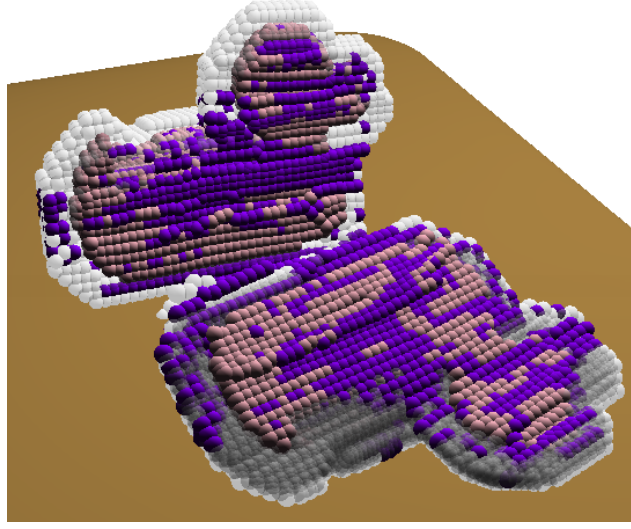


Figure 4: Particle-based damage-tracking and topological reconstruction during a cutting action. Gray and pink particles remain intact, while damaged particles (purple) mark the cut interface. From these damage signals we build an implicit SDF, extract the zero-isosurface via Marching Cubes, and then cluster connected components to recover discrete object segments.

where $J_p = \det(\mathbf{F}_p)$, ϵ_c, ϵ_s denote critical compression and stretch values, and m is a material-specific sensitivity exponent.

In practice, we select the compression and stretch thresholds ϵ_c and ϵ_s so that they correspond directly to the maximum volumetric change a particle can undergo before being flagged as damaged. Recall from Eq. (1) that

$$J_p = \det(\mathbf{F}_p)$$

measures the local volume ratio of particle p . For example, setting

$$\epsilon_c = 0.025, \quad \epsilon_s = 0.01$$

means that a particle whose volume has decreased by more than 2.5% ($J_p \leq (1 - 0.025)^m$) or increased by more than 1% ($J_p \geq (1 + 0.01)^m$) is classified as damaged. These values are chosen to be small enough to detect the onset of fracture in brittle materials, yet large enough to avoid false positives under normal elastic deformation.

The exponent m modulates how sharply damage accumulates once J_p deviates from unity: a larger m yields an abrupt transition from “healthy” to “damaged,” whereas a smaller m produces a more gradual damage accumulation. We calibrate ϵ_c , ϵ_s , and m using simulated uniaxial compression and tension tests: a small block is deformed at a constant strain rate, we record its volumetric Jacobian history, and then choose the smallest thresholds that cleanly separate reversible (elastic) behavior from irreversible damage.

For more ductile or plastic materials, one might increase the stretch threshold (e.g. $\epsilon_s \approx 0.05$) and use a lower sensitivity exponent (e.g. $m = 1$ or 2) to model gradual yielding. Conversely, for near-brittle media a tighter compression bound (e.g. $\epsilon_c \approx 0.01$) and higher exponent (e.g. $m \geq 4$) better capture sudden fracture. By anchoring these parameters to physically measurable strain limits, our framework remains both interpretable and readily tunable across a wide range of material behaviors.

For more complex yielding materials, such as those modeled with von Mises plasticity, damage occurs when the yielding stress σ_y falls below zero after particle softening:

$$\sigma_y^{(t+1)} = \sigma_y^{(t)} - \gamma \|\Delta \epsilon_p\|, \quad \text{with} \quad \text{Damaged}_p = \begin{cases} 1 & \text{if } \sigma_y^{(t+1)} \leq 0, \\ 0 & \text{otherwise,} \end{cases} \quad (2)$$

where γ is the particle softening coefficient and $\Delta \epsilon_p$ represents incremental plastic strain.

288 Once particles become damaged, their mechanical properties are altered irreversibly, enabling
 289 consistent tracking of the damaged state throughout subsequent simulation steps.

290 To reconstruct and evaluate object topology following a cut, we perform explicit topological segmen-
 291 tation leveraging the particle damage signals and knife trajectory data. We track the knife trajectory
 292 precisely during interactions, ensuring a clean spatial separation between object segments. After
 293 the completion of a cut action, defined by a stationary knife and changed particle damage states, we
 294 reconstruct the object’s surface mesh from particle data through an implicit surface representation
 295 using Signed Distance Fields (SDFs) and the Marching Cubes algorithm.

296 Formally, given the set of particles \mathcal{P} and knife trajectory \mathcal{T} , we first construct an implicit SDF
 297 representation $SDF(\mathbf{x})$ of the object’s spatial domain:

$$SDF(\mathbf{x}) = \min_{p \in \mathcal{P}} \|\mathbf{x} - \mathbf{x}_p\| - r_p, \quad (3)$$

298 where \mathbf{x}_p is the position of particle p , and r_p is its effective influence radius. Critically, we incorporate
 299 the knife trajectory \mathcal{T} by explicitly marking trajectory regions within the SDF, enforcing a strict gap
 300 devoid of particles. This guarantees that no particle fills the knife’s trajectory space, creating a clear
 301 geometric boundary between separated pieces.

302 We perform surface reconstruction by extracting the zero-isosurface of the computed SDF using
 303 a GPU-accelerated Marching Cubes implementation provided by Warp [?]. Specifically, we first
 304 define a discrete volumetric grid around particle positions, evaluate the SDF values, and then generate
 305 vertices and faces for the surface mesh:

$$\text{vertices, faces} = \text{MarchingCubes}(SDF(\mathbf{x}) = 0). \quad (4)$$

306 This mesh reconstruction process ensures explicit representation of separated object components.
 307 After extracting mesh segments, we apply Laplacian smoothing to refine surface quality:

$$\mathbf{v}_i \leftarrow \mathbf{v}_i + \alpha \sum_{j \in \mathcal{N}(i)} (\mathbf{v}_j - \mathbf{v}_i), \quad (5)$$

308 where \mathbf{v}_i is a vertex in the mesh, $\mathcal{N}(i)$ are neighboring vertices, and α is a smoothing factor.

309 Subsequently, we cluster these reconstructed mesh segments into distinct topological pieces by
 310 performing connected component analysis on the extracted mesh faces. Each particle is then assigned
 311 a cluster identity by evaluating its spatial relationship to these mesh segments. Given particle positions
 312 \mathbf{x}_p and cluster meshes, we perform the following SDF-based assignment:

$$\text{Cluster}_p = \begin{cases} i & \text{if } SDF_i(\mathbf{x}_p) < \tau, \\ -1 & \text{otherwise,} \end{cases} \quad (6)$$

313 where SDF_i is the SDF for cluster mesh i , and τ is a proximity threshold.

314 After assigning new cluster identities, we update the global object topology by comparing current
 315 particle clusters to previously stored cluster identifiers. New cluster IDs are assigned whenever
 316 new separate segments emerge, thus preserving consistent object segmentation across simulation
 317 timesteps.

318 This methodical combination of particle damage tracking, implicit surface reconstruction, and mesh-
 319 based clustering provides a robust and efficient solution for detecting and managing the topological
 320 changes that occur during robotic cutting tasks. Our approach ensures accurate evaluation of cutting
 321 outcomes essential for effective policy learning and performance assessment.

322 1.7 Pose-Invariant Shape Evaluation via Spectral Analysis in Details

323 After discovering object topology, we require a robust metric to evaluate how closely the segmented
 324 fragments match the intended goal shape. Traditional point-set distances, such as Chamfer or Earth
 325 Mover’s Distance, are sensitive to pose variations and require explicit alignment. To overcome these

limitations, we introduce a spectral-based reward derived from the intrinsic geometry of shapes, inherently invariant to rigid transformations.

Given a point cloud $X = \{x_i\}_{i=1}^n$, we construct a k -nearest neighbor graph where the edge weight between points i and j is defined as

$$W_{ij} = \exp\left(-\frac{d_{ij}^2}{\sigma^2}\right), \quad (7)$$

with d_{ij} the geodesic distance between points and σ a scaling parameter. We then build the degree matrix D and the combinatorial Laplacian:

$$D_{ii} = \sum_j W_{ij}, \quad L = D - W. \quad (8)$$

Eigen-decomposition of L yields the smallest k eigenpairs:

$$L\Phi = \Lambda\Phi, \quad (9)$$

where $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_k)$ contains eigenvalues and $\Phi = [\phi_1, \dots, \phi_k]$ are the corresponding eigenvectors. These spectral descriptors capture intrinsic shape information invariant to isometries.

Given two shapes X and Y , we define their spectral distance as:

$$d_{\text{spec}}(X, Y) = \alpha \|\Lambda_X - \Lambda_Y\|_2^2 + \beta \|\Phi_X^\top \Phi_X - \Phi_Y^\top \Phi_Y\|_F^2, \quad (10)$$

where α, β balance the contributions of eigenvalue and eigenvector differences.

To transform spectral distance into a reward, we apply a piecewise linear mapping with a critical threshold τ :

$$R(d) = \begin{cases} R_{\max} - \gamma d, & d \leq \tau, \\ R_{\max} - \gamma \tau - \delta(d - \tau), & d > \tau, \end{cases} \quad (11)$$

where γ and δ are decay rates ($\delta > \gamma$), ensuring a gradual reward decline for small errors and sharper penalties for large deviations.

For objects segmented into multiple fragments, we compute the spectral distance and associated reward for each fragment pair individually, then accumulate the total reward as:

$$R_{\text{total}} = \kappa \sum_i R(d_i), \quad (12)$$

with κ a global scaling factor ensuring consistent reward magnitudes across tasks.

This formulation provides an efficient, continuous, and pose-invariant reward signal for evaluating and planning goal-conditioned robotic cutting actions.

In order to evaluate the efficient and pose-invariant property, we conduct the following experiments:

Parameter Definitions and Pose-Invariance Observations – Our spectral-reward pipeline relies on three key configuration parameters and a fixed cutting trajectory. Specifically:

- `num_point` ($= 512$) determines how many points are uniformly sampled from the fragment’s point cloud to serve as vertices in the spectral graph. A larger `num_point` yields finer geometric resolution but increases graph-construction cost, whereas a smaller value speeds up computation at the risk of losing subtle shape details.
- k ($= 30$) defines the neighborhood size in the k -nearest-neighbor graph: each sampled point connects to its k closest neighbors (by Euclidean distance). This choice balances local connectivity (capturing fine features) against spectral stability (avoiding overly dense graphs that blur intrinsic geometry).
- *Trajectory* refers to the fixed sequence of knife motions (the “default” cutting action) applied to all shapes. At each discrete timestep along this trajectory, we recompute the fragment’s spectral descriptor and record the resulting reward value.

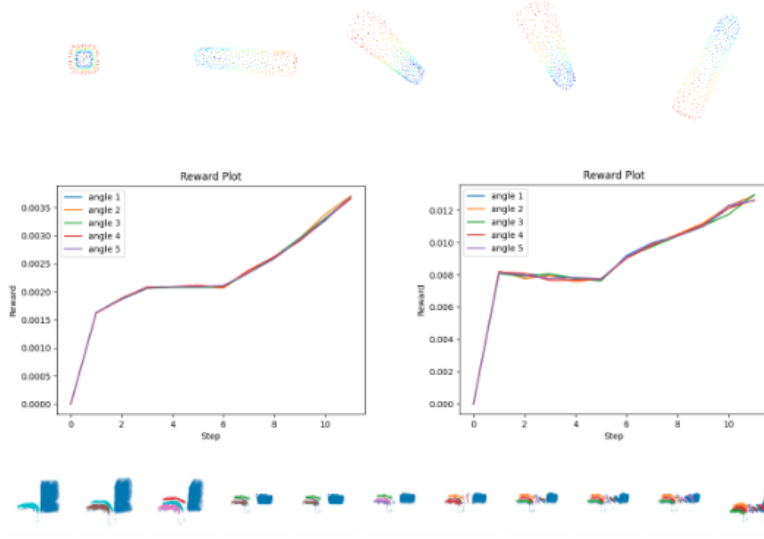


Figure 5: Pose-Invariance Evaluation of the Spectral Reward. **Top:** Five ideal stick fragments ($5 \times 5 \times 32$) rotated through distinct angles. **Middle:** Reward-versus-cut-step curves for each rotation, computed using dense particles (left) and mesh-surface particles (right) with `num_point=512` and $k = 30$. **Bottom:** Corresponding fragment point-cloud sequences along the default cutting trajectory. The nearly identical curves across all rotations confirm that our spectral reward is invariant to rigid transformations of the ideal shape.

In our experiment shown in Figure 5, we generate five identical “stick” fragments ($5 \times 5 \times 32$) and rotate each by a distinct yaw angle before sampling and evaluation. We then compute two sets of reward-vs.-step curves—one using all dense particles, the other using only mesh-surface particles. As shown in the overlaid plots, all five curves coincide almost exactly for both representations. This confirms that, under our sampling density (`num_point=512`) and graph connectivity ($k = 30$), the spectral reward is effectively invariant to rigid rotations of the ideal shape.

Multi-step Spectral-Loss Evaluation – To verify that our spectral reward correctly identifies the intended fragment at each cutting step, we conducted a pairwise matching experiment on tele-operated MPPI data. After segmenting the object into its constituent pieces, we compute the spectral distance $d_{\text{spec}}(X_{s_i}, X_{p_j})$ between the shape at cutting state s_i and each fragment p_j , for $i, j = 0, \dots, 4$. Figure 6 visualizes the resulting 4×5 “spectral-loss” matrix: rows index the initial and four successive cut states, columns index the five fragment geometries (shown below). Darker entries (lower values) indicate stronger intrinsic shape similarity. Notice that for every state s_i , the piece with low spectral distance in previous state still remain in small value despite the fact that the piece might fall down due to gravity and there is a new piece with low spectral distance appears. This demonstrates both the pose-invariance and discriminative power of our spectral descriptor, and justifies its use for automatic labeling and qualitative evaluation of MPPI rollouts.

1.8 MPPI for Data Generation

To generate high-quality demonstrations for policy learning we adopt the *Model Predictive Path Integral* (MPPI) sampling-based optimal control scheme [1]. MPPI maintains a horizon- H open-loop action sequence $\mathbf{U}_{0:H-1} \in \mathbb{R}^{H \times m}$ (with $m = 6$ for the knife pose) and repeatedly refines it by importance-sampling noisy roll-outs in the deformable-body simulator. For each iteration we draw K perturbations $\{\epsilon^{(k)}\}_{k=1}^K \sim \mathcal{N}(0, \Sigma)$ and integrate the dynamics forward, harvesting the terminal fragments $X^{(k)}$ produced by each sampled cut. The return of a roll-out is the negative spectral discrepancy

$$J^{(k)} = -R_{\text{spec}}(X^{(k)}, g) - \sum_{t=0}^{H-1} C(o_t^{(k)}, a_t^{(k)}),$$

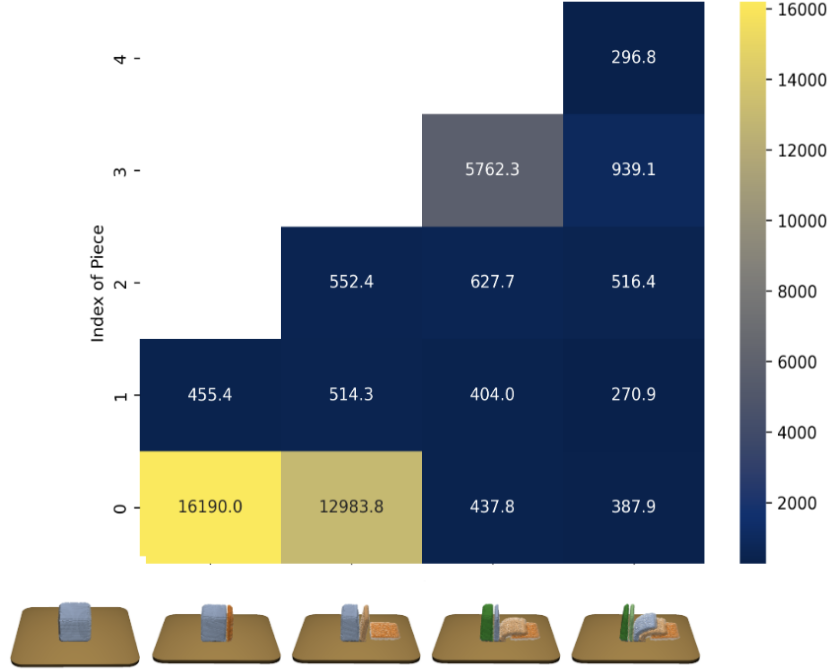


Figure 6: Pairwise spectral-loss matrix between each cutting state (rows) and each piece within the state (columns). Lower values indicate closer intrinsic geometry matches, enabling qualitative MPPI data collection.

where R_{spec} is the pose-invariant spectral reward defined previously, and $C(\cdot)$ is a small quadratic control penalty that biases the knife towards smooth, energy-efficient motions.

Given the roll-out returns we form importance weights $w_k = \exp(-J^{(k)}/\lambda)$ with temperature $\lambda = 1/\kappa$ and update the mean action sequence in closed form:

$$\mathbf{U}_{0:H-1} \leftarrow \mathbf{U}_{0:H-1} + \frac{\sum_{k=1}^K w_k \boldsymbol{\epsilon}^{(k)}}{\sum_{k=1}^K w_k}.$$

After n refinement iterations the first action \mathbf{U}_0 is executed on the robot; the horizon is then shifted and replanning continues until the episode terminates.

During each MPPI episode we log the full state-action-reward tuples $\{o_t, a_t, R_{\text{spec}}(X_t, g)\}_{t=0}^T$ as well as intermediate mesh reconstructions and knife trajectories. We repeat the procedure for a diverse set of goal shapes and material parameters, producing $\sim 6,000$ labelled cuts that serve as expert demonstrations for subsequent behaviour-cloning.

To complement the automatically generated dataset we developed a light-weight tele-operation interface in which the operator steers the knife with a 6-DoF mouse while discrete cutting commands are triggered via keyboard hot-keys. All tele-operated actions are replayed in simulation to obtain exact particle damage labels and the same spectral reward used by MPPI, ensuring perfect consistency between human and algorithmic demonstrations.

Together, MPPI planning guided by spectral rewards and targeted tele-operation yield a rich, high-fidelity dataset that underpins robust policy learning across the wide range of scenarios posed by *TopoCut*.

1.9 Policy Learning in Details

1.9.1 Dynamics-Informed Perception Module

To efficiently encode dense deformable object states for policy learning, we develop a dynamic topology prediction model, denoted as \mathcal{F} , which predicts the future topological state conditioned on

the current topological state and action. Formally, we model the dynamics as:

$$\mathcal{F} : (\text{topo}_t, a_t) \mapsto \text{topo}_{t+1}, \quad (13)$$

where both topo_t and a_t are represented as point clouds:

$$\text{topo}_t = \{(\mathbf{x}_i, \mathbf{f}_i)\}_{i=1}^N, \quad a_t = \{(\mathbf{x}_j, \mathbf{g}_j)\}_{j=1}^N, \quad (14)$$

with $\mathbf{x}_i, \mathbf{x}_j \in \mathbb{R}^3$ denoting 3D coordinates and $\mathbf{f}_i, \mathbf{g}_j$ denoting associated point features: cluster labels for topology and binary segmentation masks for action, respectively.

Within \mathcal{F} , we introduce a perception encoder Φ that processes point clouds by jointly embedding the 3D coordinates and their corresponding features:

$$\mathbf{z}_t = \Phi(\mathbf{X}_t, \mathbf{F}_t), \quad (15)$$

where \mathbf{X}_t and \mathbf{F}_t are the stacked xyz coordinates and feature vectors, respectively. Φ produces a point-wise embedding \mathbf{z}_t that captures spatial geometry, topological structure, and action intent. This embedding is later reused for downstream goal-conditioned policy learning.

Topological Representation. – For topo_t , the point features \mathbf{f}_i encode the cluster membership obtained from our particle-based damage tracking, represented as a one-hot vector over a maximum of 32 clusters. For a_t , the point features \mathbf{g}_j are binary one-hot vectors indicating whether the point lies on the cutting surface (cut = 1, not cut = 0).

Learning Objective and Intuition. – The dynamic topology model \mathcal{F} is trained to predict how the object’s topology evolves under a given action. Specifically, given (topo_t, a_t) , it predicts the resulting topological configuration topo_{t+1} . By doing so, Φ learns rich, actionable particle-wise embeddings that encode both geometric and topological transformations, enabling robust downstream manipulation.

Preprocessing and Graph Construction. – Assuming access only to noisy surface observations, we reconstruct the object’s mesh using Marching Cubes (Warp [?]), sample interior points to form volumetric point clouds, and downsample them using Farthest Point Sampling (FPS). A topology-aware graph is then constructed, where edges connect nearest neighbors within the same cluster, promoting learning of local topological consistency.

Model Architecture. – The perception encoder Φ consists of two stages: (i) a Graph Convolutional Network (GCN) to extract local features from the topology-aware graph, and (ii) a Graph Transformer that globally refines the features via self-attention. The output is a set of point-wise embeddings encoding spatial and topological context.

Training Loss. – We jointly optimize Φ with two complementary objectives: geometric consistency and topological structure prediction. The overall loss is defined as:

$$\mathcal{L} = \lambda_{\text{pos}} \mathcal{L}_{\text{pos}} + \lambda_{\text{topo}} \mathcal{L}_{\text{topo}}, \quad (16)$$

where λ_{pos} and λ_{topo} are balancing weights.

For geometric consistency, we minimize the Chamfer Distance (CD), Earth Mover’s Distance (EMD), and Hausdorff Distance (HD) between predicted and ground-truth point clouds:

$$\mathcal{L}_{\text{pos}} = \text{CD}(\hat{X}, X) + \text{EMD}(\hat{X}, X) + \text{HD}(\hat{X}, X) \quad (17)$$

where \hat{X} and X are predicted and ground-truth point sets.

For topological structure prediction, we treat clustering as a contrastive learning task, where only the relative membership between points matters. A Hungarian matching-based loss is applied:

$$\mathcal{L}_{\text{topo}} = \min_{\pi \in \mathcal{P}} \sum_{i,j} \ell(\pi(c_i), \hat{c}_j), \quad (18)$$

where \mathcal{P} is the set of bipartite matchings, c_i and \hat{c}_j are the ground-truth and predicted cluster assignments, and $\ell(\cdot, \cdot)$ is a binary cross-entropy loss.

444 *Training Data.* – The dynamic topology model \mathcal{F} is pretrained using state-action-state triplets
 445 $(\text{topo}_t, a_t, \text{topo}_{t+1})$ collected from the MPPI-based demonstration generation (Section ??). This
 446 rich dataset covers diverse materials, object geometries, and cutting scenarios, enabling the learned
 447 embeddings to generalize effectively across a wide range of deformable manipulation tasks.

448 1.9.2 Particle-based Score-Entropy Discrete Diffusion Policy

449 With the pretrained perception encoder Φ that captures geometric and topological information, we
 450 now train a goal-conditioned behavior cloning (BC) policy for robotic cutting. The policy operates
 451 over perception embeddings, taking as input the current object topology and the desired goal shape,
 452 and predicting the next cutting action.

453 Given a point cloud observation $o_t = (\mathbf{X}_t, \mathbf{F}_t)$ representing the current object state, and a goal point
 454 cloud $g = (\mathbf{X}_g, \mathbf{F}_g)$, we first obtain their embeddings:

$$\mathbf{z}_t = \Phi(\mathbf{X}_t, \mathbf{F}_t), \quad \mathbf{z}_g = \Phi(\mathbf{X}_g, \mathbf{F}_g), \quad (19)$$

455 where \mathbf{z}_t and \mathbf{z}_g are point-wise embeddings encoding geometry and topology.

456 The cutting action is represented as a binary segmentation $\mathbf{a}_t \in \{0, 1\}^N$, where each point is classified
 457 as cut (1) or not cut (0).

458 *Policy Model: Conditional Score-Based Discrete Diffusion.* – Inspired by Score Entropy Discrete
 459 Diffusion (SEDD) [2], we formulate action prediction as a conditional discrete denoising diffusion
 460 process over point-wise binary labels.

461 In the forward process, the clean action labels \mathbf{a}_t^* are progressively corrupted into noisy labels $\tilde{\mathbf{a}}_t$
 462 through a multinomial noise distribution:

$$q_t(\tilde{\mathbf{a}}_t | \mathbf{a}_t^*) = \text{Multinomial}(\tilde{\mathbf{a}}_t | \mathbf{p}_t(\mathbf{a}_t^*)), \quad (20)$$

463 where $\mathbf{p}_t(\mathbf{a}_t^*)$ denotes the noise schedule at timestep t .

464 The policy network s_θ is trained to predict the score function:

$$s_\theta(\tilde{\mathbf{a}}_t, t, \mathbf{z}_t, \mathbf{z}_g) \approx \nabla_{\tilde{\mathbf{a}}_t} \log q_t(\mathbf{a}_t^* | \tilde{\mathbf{a}}_t), \quad (21)$$

465 where s_θ outputs the gradient of the log-posterior of the clean action given the noised action,
 466 conditioned on the current and goal embeddings.

467 This formulation enables the policy to iteratively denoise $\tilde{\mathbf{a}}_t$ toward recovering the target cutting
 468 action.

469 *Training Objective.* – The BC policy is trained by minimizing the Denoising Score Entropy (DSE)
 470 loss across all diffusion steps:

$$\mathcal{L}_{\text{BC}}(\theta) = \mathbb{E}_{(o_t, g, \mathbf{a}_t^*) \sim D} \left[\sum_{t=1}^T \mathcal{L}_{\text{DSE}}(\tilde{\mathbf{a}}_t, \mathbf{a}_t^*; \theta) \right], \quad (22)$$

471 where the per-step DSE loss is defined as:

$$\mathcal{L}_{\text{DSE}}(\tilde{\mathbf{a}}_t, \mathbf{a}_t^*; \theta) = \mathbb{E} \left[\|s_\theta(\tilde{\mathbf{a}}_t, t, \mathbf{z}_t, \mathbf{z}_g) - \nabla_{\tilde{\mathbf{a}}_t} \log q_t(\mathbf{a}_t^* | \tilde{\mathbf{a}}_t)\|_2^2 \right]. \quad (23)$$

472 *Action Reconstruction.* – After completing the denoising sampling process, we obtain a binary
 473 segmentation over the points indicating which regions should be cut. From the points classified as
 474 cut, we fit a cutting plane, thereby reconstructing the knife’s pose and producing the final action a_t .

475 *1. Inputs, Goals, and Perception Embeddings* – At each decision step t , the robot observes:

$$o_t = (\mathbf{X}_t, \mathbf{F}_t),$$

476 where:

477 • $\mathbf{X}_t = [\mathbf{x}_{t,1}, \dots, \mathbf{x}_{t,N}]^\top \in \mathbb{R}^{N \times 3}$ are the 3D coordinates of N particles representing the object.

478 • $\mathbf{F}_t = [\mathbf{f}_{t,1}, \dots, \mathbf{f}_{t,N}]^\top \in \mathbb{R}^{N \times f}$ are associated per-point features, such as normals, damage
 479 indicators, or material properties.

480 A desired goal shape is provided as another point-cloud:

$$g = (\mathbf{X}_g, \mathbf{F}_g) \in \mathbb{R}^{N \times 3} \times \mathbb{R}^{N \times f}.$$

481 Both o_t and g are passed through a pretrained perception encoder

$$\Phi : (\mathbb{R}^{N \times 3}, \mathbb{R}^{N \times f}) \longrightarrow \mathbb{R}^{N \times d},$$

482 which uses graph-based convolutions followed by a small MLP to produce point-wise latent embed-
 483 dings:

$$\mathbf{z}_t = \Phi(\mathbf{X}_t, \mathbf{F}_t), \quad \mathbf{z}_g = \Phi(\mathbf{X}_g, \mathbf{F}_g),$$

484 where d is the embedding dimension. Intuitively, \mathbf{z}_t captures the current object’s geometric and
 485 topological state, while \mathbf{z}_g encodes the desired target shape.

486 *2. Action Representation and Forward Noising* – We represent the cutting action at time t as a binary
 487 mask

$$\mathbf{a}_t^* = [a_{t,1}^*, \dots, a_{t,N}^*]^\top \in \{0, 1\}^N,$$

488 where $a_{t,i}^* = 1$ indicates that particle i should be cut. To train our policy using denoising diffusion, we
 489 first define a forward noising process that corrupts \mathbf{a}_t^* into progressively noisier versions $\tilde{\mathbf{a}}_1, \dots, \tilde{\mathbf{a}}_T$.
 490 Starting from $\tilde{\mathbf{a}}_0 = \mathbf{a}_t^*$, each step applies a small random flip:

$$q_t(\tilde{\mathbf{a}}_t \mid \tilde{\mathbf{a}}_{t-1}) = \prod_{i=1}^N \left[(1 - \beta_t) \delta(\tilde{a}_{t,i} = \tilde{a}_{t-1,i}) + \beta_t \delta(\tilde{a}_{t,i} \neq \tilde{a}_{t-1,i}) \right],$$

491 where $\beta_t \in [0, 1]$ is a noise-rate schedule (e.g. $\beta_t = t/T$). After T steps, $\tilde{\mathbf{a}}_T$ is nearly uniform
 492 random.

493 *3. Score Network and Reverse Diffusion* – We train a neural network $s_\theta(\tilde{\mathbf{a}}_t, t, \mathbf{z}_t, \mathbf{z}_g)$ to predict the
 494 gradient of the log-probability of the true mask given the noisy mask:

$$s_\theta(\cdot) \approx \nabla_{\tilde{\mathbf{a}}_t} \log q_t(\mathbf{a}_t^* \mid \tilde{\mathbf{a}}_t).$$

495 Concretely, s_θ receives as input:

- 496 • The noisy mask $\tilde{\mathbf{a}}_t$ (embedded as scalars).
- 497 • The timestep t , encoded with sinusoidal features.
- 498 • The concatenated perception embeddings $[\mathbf{z}_t; \mathbf{z}_g] \in \mathbb{R}^{N \times 2d}$.

499 During inference, we reverse the noising chain by sampling

$$p_\theta(\mathbf{a}_{t-1} \mid \tilde{\mathbf{a}}_t, t, \mathbf{z}_t, \mathbf{z}_g) = \text{Categorical}\left(\text{softmax}(\log(1 - \beta_t, \beta_t) + s_\theta(\tilde{\mathbf{a}}_t, t, \mathbf{z}_t, \mathbf{z}_g))\right).$$

500 This step “denoises” $\tilde{\mathbf{a}}_t$ one level at a time back toward a clean mask.

501 *4. Training Objective: Denoising Score-Entropy Loss* – We optimize θ by minimizing the expected
 502 squared error between the network’s predicted score and the true score over all timesteps:

$$\mathcal{L}_{\text{BC}}(\theta) = \mathbb{E}_{(o_t, g, \mathbf{a}_t^*) \sim D} \left[\sum_{t=1}^T \underbrace{\mathbb{E}_{\tilde{\mathbf{a}}_t \sim q_t(\cdot \mid \mathbf{a}_t^*)} \| s_\theta(\tilde{\mathbf{a}}_t, t, \mathbf{z}_t, \mathbf{z}_g) - \nabla_{\tilde{\mathbf{a}}_t} \log q_t(\mathbf{a}_t^* \mid \tilde{\mathbf{a}}_t) \|_2^2}_{\mathcal{L}_{\text{DSE}}(\tilde{\mathbf{a}}_t, \mathbf{a}_t^*; t)} \right].$$

503 Because q_t uses simple bit-flip noise, the true score gradient has a closed form:

$$\nabla_{\tilde{a}_{t,i}} \log q_t(a_i^* \mid \tilde{a}_{t,i}) = \frac{\delta(\tilde{a}_{t,i} = a_i^*) - (1 - \beta_t)}{\beta_t(1 - \beta_t)}.$$

504 *5. Inference and Action Reconstruction* – At test time:

- 505 1. Initialize $\tilde{\mathbf{a}}_T$ by sampling each bit from Bernoulli(0.5).

- 506 2. For $t = T, T - 1, \dots, 1$, sample \mathbf{a}_{t-1} from $p_\theta(\mathbf{a}_{t-1} \mid \tilde{\mathbf{a}}_t, t, \mathbf{z}_t, \mathbf{z}_g)$.
507 3. The final mask $\tilde{\mathbf{a}}_0$ indicates the cut points.
508 4. Fit a planar cut by solving

$$\min_{\mathbf{n}, d} \sum_{i: \tilde{a}_{0,i}=1} (\mathbf{n}^\top \mathbf{x}_{t,i} + d)^2, \quad \|\mathbf{n}\|_2 = 1,$$

509 where \mathbf{n} is the plane normal and d its offset.

510 *Summary.* – Our BC policy systematically integrates dynamics-informed perception embeddings
511 and conditional score-based discrete diffusion modeling. By treating cutting action prediction as a
512 goal-conditioned denoising process over discrete labels, our method achieves robust and generalizable
513 cutting behavior across diverse object shapes and cutting goals.

514 PDDP naturally handles the combinatorial nature of cutting (multimodal mask distributions) and
515 provides smooth trade-offs between sample quality and runtime via T . The closed-form score
516 supervision ensures stable training, and the iterative denoising generates crisp, coherent cutting
517 actions that generalize across object shapes and materials.

518 Overall, our method systematically combines efficient topology tracking, pose-invariant shape
519 evaluation, dynamics-informed perception, and supervised policy learning to enable effective robotic
520 cutting in general scenarios.

521 References

- 522 [1] G. Williams, P. Wagener, B. Goldfain, P. Drews, J. Rehg, and E. Theodorou. Information-
523 theoretic model predictive path integral control: Theory and applications to robotics. In *Proceed-*
524 *ings of Robotics: Science and Systems (RSS)*, 2017.
- 525 [2] A. Lou, C. Meng, and S. Ermon. Discrete diffusion modeling by estimating the ratios of the data
526 distribution, 2024. URL <https://arxiv.org/abs/2310.16834>.