

373 6 Supplementary Materials

374	6.1 Perception	13
375	6.1.1 Data Collection	13
376	6.1.2 Data Preprocessing	13
377	6.2 Dynamics Model	14
378	6.2.1 Graph Building	14
379	6.2.2 Model Training	14
380	6.3 Closed-loop Control	14
381	6.3.1 Action Space	14
382	6.3.2 Multi-bin Classification	15
383	6.4 Tool Design	15
384	6.5 Tool-Switching Setup	17
385	6.6 Comparison with Human Subjects on Dumpling-making	17
386	6.7 Tool Classification	17
387	6.8 Human Evaluation of Alphabet Letters	17
388	6.9 Reinforcement Learning Baseline	18

389 6.1 Perception

390 6.1.1 Data Collection

391 To train the GNNs, we collect around 20 minutes of data for each tool using a behavior policy that
392 randomly samples parameters within a predefined action space. In each episode, the robot applies
393 multiple action sequences to the soft object, and after the current episode, a human resets the dough’s
394 shape into another random form.

395 The data collected for each tool are as follows: (1) Asymmetric gripper / two-rod symmetric
396 gripper / two-plane symmetric gripper: 60 episodes with five sequences per episode; (2) Circle
397 press / square press / circle punch / square punch: 90 episodes with three sequences per episode; (3)
398 Large roller / small roller: 80 episodes with three sequences per episode.

399 We collect point clouds before and after executing each sequence to train the tool classifier. How-
400 ever, we augment the training data by including any pair of these point clouds, not just consecutive
401 pairs. For tools that don’t require a GNN-based dynamics model, we execute a pre-coded dumpling-
402 making pipeline ten times and add point clouds captured before and after using each tool in the
403 pipeline to our dataset. Note that the majority of tool selection data is a direct reuse of the dynamics
404 data collected during the training of the dynamics model. This approach efficiently collects real-
405 world tool selection data without needing extra exploration. We record the data collection process
406 in the fourth supplementary video.

407 6.1.2 Data Preprocessing

408 When building the dataset to train the dynamics model, aside from optimizing the sample quality at
409 each time frame, we also want to leverage the continuity of the video data. Therefore, we introduce
410 simple geometric heuristics into the physical environment for better frame consistency. First, if the
411 operating tool is not in contact with the convex hull of the object point cloud, we use the same
412 sampled particles from the previous frame. This also applies when the tool moves away from the
413 object. Additionally, we subsample the original videos to ensure that each video in the dataset has
414 the same number of frames (16 frames in practice).

415 6.2 Dynamics Model

416 6.2.1 Graph Building

417 When building the graph, the edges between the dough particles are constructed by finding the
418 nearest neighbors of each particle within a fixed radius (in practice, 0.1 cm). The edges between
419 the tool and dough particles are computed slightly differently. Instead of simply connecting to all
420 the neighbors within the threshold, we limit the number of undirected edges between tool particles
421 and the dough particles to at most four per tool particle to cut off the redundant edges in the graph
422 neural network. Since all the node and edge features in the GNN are encoded in each particle’s local
423 neighborhood, our GNN is naturally translation-invariant and therefore can accurately predict the
424 movement of the dough regardless of its absolute location in the world frame.

425 6.2.2 Model Training

426 We train the model with temporal abstraction to enhance performance and inference speed. For
427 example, when $t = 0$, we train the model to predict the state of the dough at $t = 3$ directly instead
428 of $t = 1$. This shortens the horizon, eases the task, and improves our model’s inference speed by
429 decreasing the number of forward passes needed for a full action sequence.

430 6.3 Closed-loop Control

431 6.3.1 Action Space

432 We classify the tools into a few categories based on $|\mathcal{A}|$, the dimension of their corresponding action
433 space. We visualize action spaces for gripping, pressing, and rolling in Figure 3 (B).

434 A) Nine tools that have an action space with $|\mathcal{A}| \geq 3$:

- 435 1) Asymmetric gripper / two-rod symmetric gripper / two-plane symmetric gripper:
436 $\{r, \theta, d\}$, where r is the distance between the midpoint of the line segment connecting
437 the centers of mass of the gripper’s two fingers and the center of the target object, θ
438 is the robot gripper’s rotation about the (vertical) axis, and d is the minimal distance
439 between the gripper’s two fingers during this pinch.
- 440 2) Large roller / small roller / square press / square punch: $\{x, y, z, \theta\}$, where $\{x, y, z\}$
441 is the bounded location indicating the center of the action, and θ is the robot gripper’s
442 rotation about the vertical axis. In the case of rollers, the rolling distance is fixed and
443 therefore not included in the action space.
- 444 3) Circle press / circle punch: $\{x, y, z\}$, where $\{x, y, z\}$ is the bounded location indicat-
445 ing the center of the action. The robot gripper’s rotation is unnecessary because the
446 tool’s bottom surface is a circle.

447 B) Five tools that have an action space with $|\mathcal{A}| = 2$:

- 448 1) Knife / circle cutter / pusher / skin spatula / filling spatula: $\{x, y\}$, where $\{x, y\}$ is the
449 bounded location indicating the center of the action on the plane. θ and z are fixed for
450 these tools to simplify the action space.

451 C) The action of the hook is precoded.

452 In category B, for all tools except the knife, we leverage the prior that the center of the dough is
453 always the optimal solution in the action space and directly compute the center from the processed
454 point cloud. In the case of the knife, we use the y coordinate of the center of the dough as the
455 solution for y (the xyz coordinate system is illustrated in Figure 4). For x , we first compute the
456 volume of the target dough and then perform a binary search with the center of the dough as the
457 starting point to find the cut position that results in the closest volume to the target volume.

458 In category C, the hook is precoded first to hook the handle of the dumpling mold, then close the
459 mold, press the mold handle to turn the dough into a dumpling shape, and finally open the mold by
460 hooking and lifting the handle.

461 The guiding principle in designing action spaces involves starting with the end-effector’s 6-DoF
 462 action space and eliminating redundant DoFs. For instance, rotations along the x and y axes are
 463 typically not required to generate a meaningful action. Hence, we opt to exclude them from the
 464 action space of the 14 tools. For grippers, we transform the Cartesian coordinate system into a polar
 465 coordinate system to simplify the search process for action parameters since corner cases in the
 466 bounded Cartesian space are usually suboptimal. Following this, we introduce tool-specific DoFs,
 467 which are determined by the tool’s geometric properties. For example, in the case of grippers, we
 468 incorporate an additional parameter, d , to represent the width between the gripper’s two fingers.

469 Our method can potentially generalize to various challenging dough manipulation tasks besides
 470 dumpling-making, such as making alphabet letter cookies (as shown in the paper), pizza, and noo-
 471 dles. A successful transfer requires the ground truth meshes of new tools and data from interacting
 472 with them. We only need 20 minutes of real-world interaction data per tool, demonstrating the ease
 473 of retraining for new tasks and tools. Although we incorporate human prior knowledge to simplify
 474 the action space for tools, it does not constrain the generalization capability since we can easily
 475 specify the action space for new tools.

476 6.3.2 Multi-bin Classification

477 We formulate the self-supervised policy training as a multi-bin classification problem inspired by
 478 previous works on 3D bounding box estimation [50, 51]. The total loss for the multi-bin classifica-
 479 tion is

$$\mathcal{L} = \sum_{i=1}^{|\mathcal{A}|} \left(\mathcal{L}_{\text{conf}}^{\mathcal{A}_i} + w \cdot \mathcal{L}_{\text{loc}}^{\mathcal{A}_i} \right), \quad (5)$$

480 where the confidence loss $\mathcal{L}_{\text{conf}}^{\mathcal{A}_i}$ is the softmax loss of the confidences of each bin for each action
 481 parameter \mathcal{A}_i , and the localization loss $\mathcal{L}_{\text{loc}}^{\mathcal{A}_i}$ is the loss that tries to minimize the difference between
 482 the estimated parameter and the ground truth parameter. For orientation estimation, we use negative
 483 cosine loss as the localization loss and force it to minimize the difference between the ground truth
 484 and all the bins that cover that value. We use the smooth L1 loss as the localization loss for action
 485 parameters not representing an orientation. During inference time, for each parameter, the bin with
 486 maximum confidence is selected, and the final output is computed by adding the estimated delta of
 487 that bin to the center of the same bin.

488 6.4 Tool Design

489 We design and 3D-print 14 tools: large roller, small roller, circle press, circle punch, square press,
 490 square punch, knife / pusher, circle cutter, two-rod symmetric gripper, asymmetric gripper, two-
 491 plane symmetric gripper, skin spatula, filling spatula, and hook. The dumpling mold is the same
 492 as real-world ones. In Figure 7, we compare our 3D-printed tools and their real-world prototypes,
 493 which are common kitchen tools for dough manipulation. The design principle of these 3D-printed
 494 tools is to mimic real-world ones as closely as possible.

495 The roller is composed of a holder and a rolling pin so that the rolling pin can rotate freely while
 496 the holder remains static. We designed both large and small rollers to accommodate different needs.
 497 We also have a set of punches and presses with square and circle shapes. The knife is a thin planar
 498 tool that can cut through objects. Similarly, the circle cutter can cut an object into a circular shape.
 499 Among the grippers, the two-rod symmetric gripper consists of two cylindrical extrusions, the asym-
 500 metric gripper consists of a cylindrical and planar part, and the two-plane symmetric gripper consists
 501 of two planar parts. The two extruding rods on each gripper insert into the corresponding holes of
 502 the two fingers of Franka’s gripper, allowing them to adhere to and move along with the fingers.
 503 A linear shaft connects the two parts of each gripper, constraining their movement to a single axis.
 504 The skin and filling spatulas have a similar design, except that their two extrusions are each spatula,
 505 so they can pick up and put down the soft object without deforming it. The hook and the dumpling
 506 mold are tools used together to mold the dough into a dumpling shape.

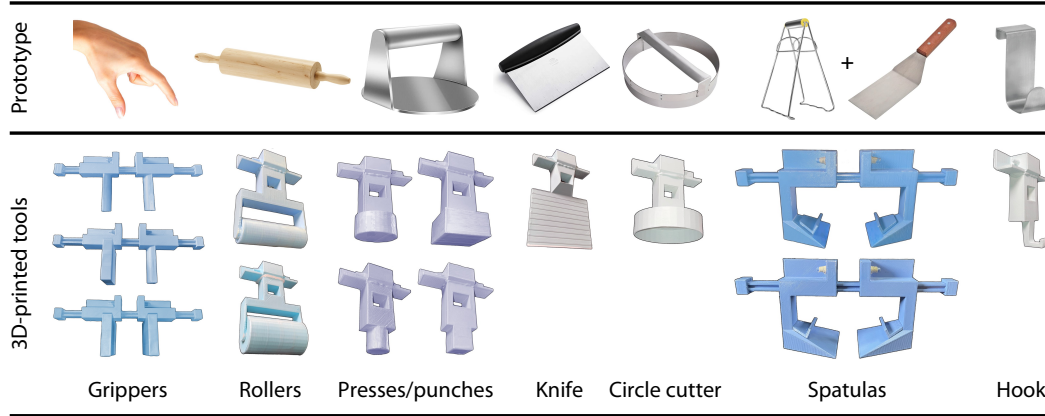


Figure 7: **Prototypes of 3D-printed tools.** We show a comparison between our 3D-printed tools and their real-world prototypes which are common kitchen tools for dough manipulation. The design principle of these 3D-printed tools is to mimic real-world ones as closely as possible. We use 3D-printed tools instead of real-world ones to allow the robot arm to acquire and manipulate the tools more easily.

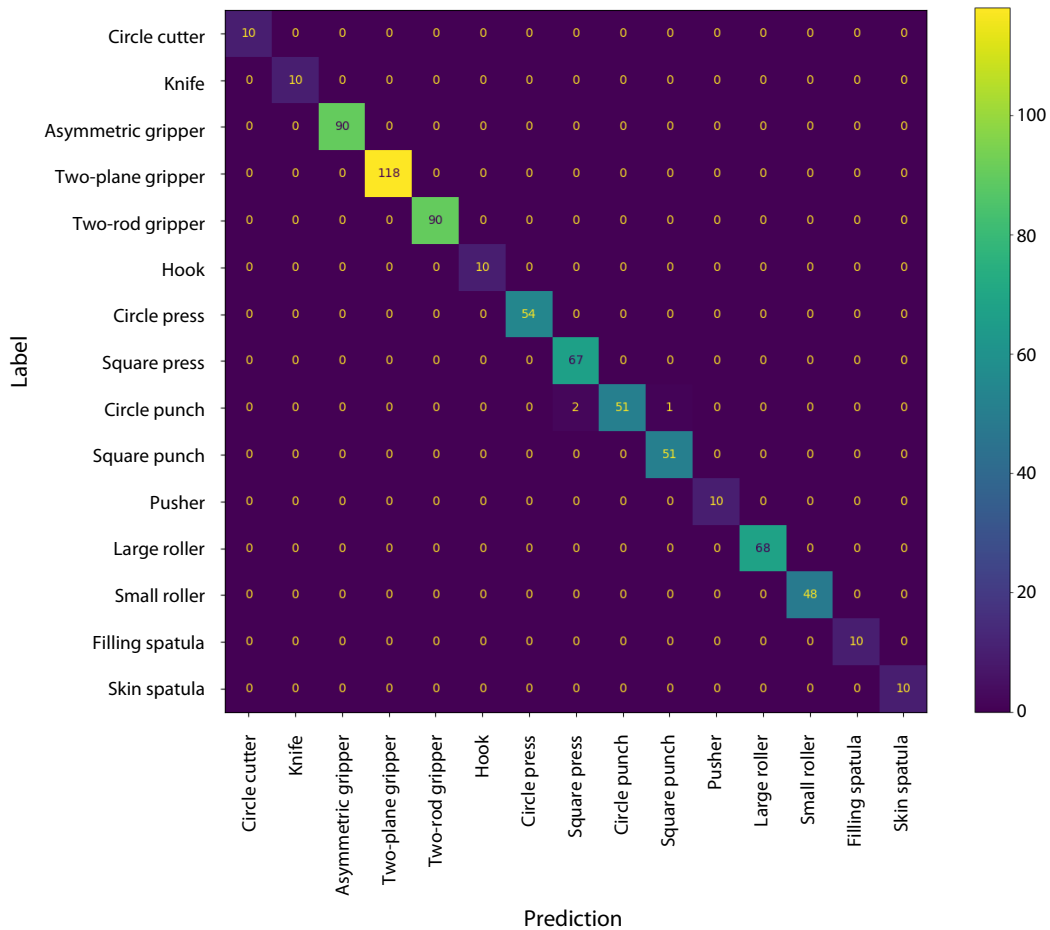


Figure 8: **Confusion matrix of the tool classifier predictions.** We show the confusion matrix of the tool classifier predictions on the test set, which is split from the training data. The tool classifier achieves an accuracy very close to 1.

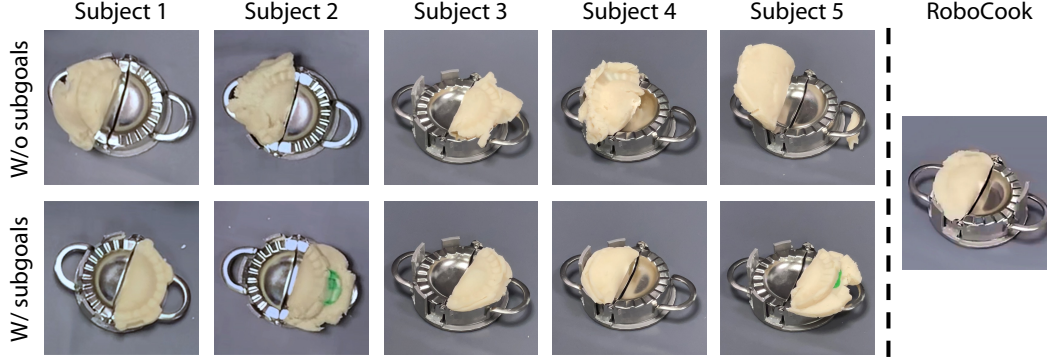


Figure 9: **Comparison with human subjects.** We show a comparison with the manipulation results of human subjects. In the first row, Human subjects devise their manipulation plan and choose tools independently. In the second row, human subjects follow a given tool sequence and subgoals.

6.5 Tool-Switching Setup

The tool-switching setup is an engineering innovation we implement in this project. We adopt two designs so that the robot can pick up, use, and put down the tools without any help from humans: (1) The connector on the top of each tool attaches firmly to the Franka’s gripper when it closes its fingers and also unlocks easily when the gripper reopens. (2) The tool racks on the bottom and right side of the robot table hold all the 3D-printed tools in their upright poses so that the robot can easily pick up and put down the tools. Additionally, we calibrate the tools’ world coordinates so that the robot knows where to find each tool. The supplementary videos of making dumplings show examples of how the robot switches tools.

6.6 Comparison with Human Subjects on Dumpling-making

We invited five human subjects to make dumplings with the same tools to highlight the complexity of dumpling-making. Each subject participated in two experiments: choosing tools independently and following a given tool sequence and subgoals. For a fair comparison, human subjects were not allowed to directly touch the dough with their hands or apply each tool more than five times. Before the experiments, we introduced each tool and gave them sufficient time to get familiar with the dough’s dynamics and devise their plan. We compare their best attempt among three trials to our method for each experiment. Figure 9 shows that human performance is notably worse than our method without subgoals. Performance improves with the tool sequence and subgoals but remains comparable to or worse than our method. The fifth supplementary video records the entire process.

6.7 Tool Classification

We split a test set from the training data of the tool classifier and show the confusion matrix of the tool classifier predictions in Figure 8. The instance accuracy is 0.996. We compared PointNet-based and ResNet-based architectures for the tool classification network. PointNet-base architecture generalizes better due to its ability to encode depth information. Empirically, it demonstrates greater robustness to changes in lighting, dough color, and dough transformations.

6.8 Human Evaluation of Alphabet Letters

We recognize a discrepancy between how metrics such as Chamfer Distance measure the results and how humans perceive them - these metrics are prone to local noises while humans are good at capturing the holistic features of the dough. Therefore, we invite 100 human subjects to evaluate the results. The human survey asks the question: “What alphabet letter is the robot trying to shape in the given image?” If we put all 20 images (four methods × five letters) in Question 1, there could be a predictive bias from seeing more than one image of the same letter. Therefore, we shuffle the order

539 of 20 images and split them into four groups. Each group contains one image for each letter but
540 from different methods. After averaging over five letters, we show the human perceived accuracy
541 and human ranking of the performance of these four methods in Table 1.

542 **6.9 Reinforcement Learning Baseline**

543 The RL+GNN baseline utilizes a model-based Soft Actor-Critic (SAC) with a learned GNN-based
544 dynamics model as the world model. The action space aligns with other planning methods, and the
545 state space comprises the point cloud position. The reward function is derived from the change in
546 Chamfer Distance after each grip. Training involves a discount factor of 0.99, a learning rate of
547 0.0003 with the Adam optimizer, 2-layer MLPs with 256 hidden units, and ReLU activation for both
548 policy and critic models. We initially collect 250 steps of warm-up data. The replay buffer size is
549 $1e6$, and the target smooth coefficient is 0.005. The results shown in Figure 5 and Table 1 indicate
550 that the RL baseline is noticeably worse compared to our method.