

Supplementary Material

SCONE: A Food Scooping Robot Learning Framework with Active Perception

Anonymous Author(s)

Affiliation

Address

email

1 The supplementary materials consist of:

- 2 1. Demonstration video.
- 3 2. Details of dataset.
- 4 3. Baseline approach.
- 5 4. Analysis of experiment results.
- 6 5. Details of our proposed method.
- 7 6. Visualization of attention scores.

8 **1 Demonstration video**

9 In the supplementary video, we show 1) a brief introduction of our food scooping robot learning
10 framework and 2) illustrations of scooping tasks and qualitative results.

11 **2 Details of Dataset**

12 In this section, we provide more details about the data collection and preparation process for training.

13 **2.1 Food Preparation**

14 We select a total of 12 categories of food items for our real-world experiment, with 6 categories used
15 for training and the remaining 6 categories used for testing. To simplify the complexity, we have
16 limited the differences between categories primarily to **particle size** and the **amount** of food items.

17 **Food in the Training Set.** The food categories for training contain brown rice, mung bean, soybean,
18 chocolate call, dried jujube, and cheese ball. For food items with small particles such as brown rice,
19 mung bean, and soybean, we fill the bowl to approximately 2/3 of its capacity. This threshold ensures
20 that the interaction can be carried out without the risk of spilling the food items. For food items with
21 large particles such as chocolate balls, dried jujube, and cheese balls, we select more than one piece
22 but fewer than a certain number based on their size. This is done to prevent the spoon from getting
23 stuck or breaking the food items when there is insufficient space in the bowl for the spoon to reach
24 them.

25 **Food in the Testing Set.** The food categories for testing contain sago, red bean, orange, macadamia,
26 penne, and fruit candy. In the testing set, we have designed three levels of difficulty to evaluate the
27 performance of the models. Both the **Basic** and **Extended** settings in the evaluation include sago,
28 red bean, orange, and macadamia, and the peculiar setting includes penne and fruit candy.

- 29 • **Basic Setting:** In this setting, the conditions are kept identical to the training set. For food items
30 with small particles, the bowl is filled to approximately $2/3$ of its capacity. For food items with large
31 particle sizes, more than one entity is included in the bowl.
- 32 • **Extended Setting:** In this setting, we change the combination of properties related to particle size
33 and amount of food items. We want to explore different scenarios to evaluate the performance of our
34 model under varying conditions. For food items with small particles, we fill the bowl with a smaller
35 quantity of these items. The intention is to keep the height of the food in the bowl similar to that of
36 the food items with large particles in the **Basic** setting. Conversely, we increase the amount of food
37 items with large particles.
- 38 • **Peculiar Setting:** These food items in the peculiar setting had unique features, such as different
39 shapes, colors, or textures, that are not present in the training set. By introducing these visually
40 distinct food items, we aim to challenge the model’s capacity to recognize and handle novel objects
41 effectively.

42 2.2 Manipulation Policy

43 During the data collection phase, we employ two distinct manipulation policies to ensure the suc-
44 cessful scooping of the food items. To prevent spilling, we adopt a specific policy for scooping
45 up food items with small particles. The strategy involves positioning the spoon at a shallow depth
46 under the height of the food items. By adopting this approach, the risk of spillage can be minimized
47 and the food items were securely contained within the spoon during the scooping process. When
48 dealing with food items that have large particles, we position the spoon at the lowest point within the
49 bowl and gently push the items toward the edge. This allows the food to roll into the spoon and is
50 able to successfully scoop up without any spillage. Figure 1 shows the end-effector positions during
51 manipulation in the training dataset.

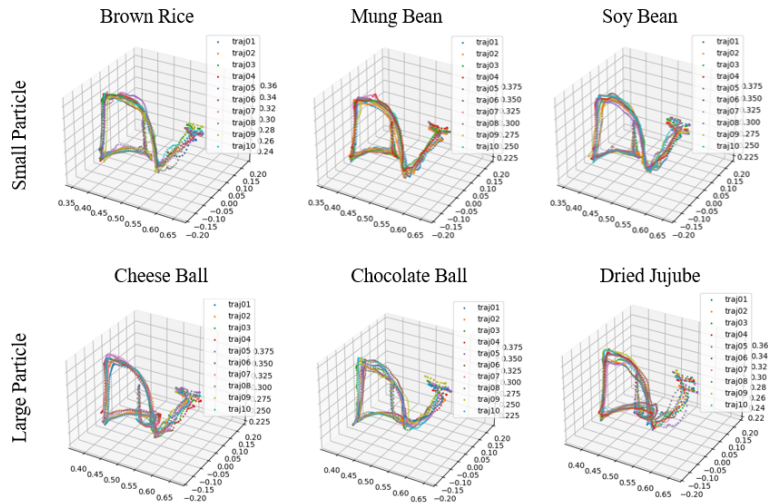


Figure 1: Visualization of End-effector Position During Manipulation.

52 2.3 Implementing Detail of Overall Scooping Task

53 The task is divided into three stages: interacting, scooping, and transferring. Among these stages,
54 only the scooping stage requires learning. The overall task follows a predefined procedure. The
55 interaction data is collected by replaying the recorded end-effector trajectory. The model then be-
56 gins predicting the scooping stage. Once the scooping stage is completed, the transferring stage is
57 initiated. The trajectory for the transferring stage is obtained by averaging all training data. Consid-
58 ering that the ending position of the scooping stage for every testing food tends to be relatively close
59 to initial position of the transferring stage, the Panda arm then proceeds to move to the first point

60 using simple motion planning. Once the transition is complete, the rest of the trajectory is replayed
61 accordingly.

62 3 Baseline Approach

63 This section provides detailed explanations of the baselines in comparison with SCONE.

64 3.1 BC-based Method

65 When adopting Learning from Demonstration (LfD), behavior cloning (BC) is a straightforward
66 approach for direct learning from both observation and action. We conducted the task using 4
67 different BC models, both with and without the inclusion of interacting data.

68 **BC one stage.** In BC (one stage), we consider scooping and transferring as a single stage. This
69 implies that the BC (**one stage**) model is required to learn the long-horizon task only through obser-
70 vation, without any additional input.

71 **BC without interaction.** The most basic method that learns the scooping task through observations.

72 **BC with food ID.** To implement BC with food ID, we applied the trained classification model to the
73 interaction data. The predicted food ID was then conditioned as a one-hot vector and concatenated
74 with the observation z_{gi} , serving as the input for the BC model.

75 **BC with interaction.** In the case of BC with interaction, the input for the model is the concatenation
76 of latent of observation and interaction data.

77 3.2 Template Policy

78 One of the approaches we explore in utilizing interaction data is the selection of trajectories from
79 the pre-interacting data based on their similarity to the food categories in the training dataset. The
80 template trajectory is obtained by averaging the sequences of end-effector pose of 6 different food
81 categories respectively.

82 **Rand. Template.** Random policy is selected arbitrarily.

83 **Classified Template.** The template policy is selected accordingly based on the predicted category
84 by the classifier.

85 3.3 Dynamical System Model

86 **MTRNN.** We utilized the multiple time scale recurrent neural network by [1] to update the initial
87 parameter Cs_0 using the interaction data in the testing stage.

88 4 Analysis of Experiment Results

89 4.1 Failure Cases

90 During our real-world evaluation, the failure cases observed included instances of spilling (SP),
91 insufficient food on the spoon (IF), failed attempts to scoop (FA), collisions (CO), and others (OT).
92 We select the orange, sago, and penne as examples, and their results can be found in [Table 1](#) and
93 [Table 2](#).

94 **Spilling (SP).** Spilling happens at the end of the scooping process, especially when there is an
95 excessive amount of food on the spoon. Consequently, during the stages of transferring or when
96 withdrawing the spoon from the bowl, the excess food spills out.

97 **Insufficient food (IF).** The failure case of insufficient food occurs when the amount of food on the
98 spoon is unable to cover at least one-third of its surface area. This is frequently attributed to the
99 spoon not being inserted deep enough into the bowl to effectively reach the desired food portion.

	Basic - Orange					
	SP	IF	FA	CO	OT	Failed
BC (one stage)	0	0	5	3	1	9
BC (w/o interact)	2	0	3	3	0	8
BC (w/ food id)	1	0	2	1	1	5
BC (w/ interact)	1	0	2	0	0	0
Rand. Template	0	0	5	0	0	5
Classified Template	0	0	0	0	0	0
MTRNN [1]	0	0	0	0	0	0
SCONE (Ours)	0	0	0	0	0	0

(a) Basic - Orange

	Basic - Sago					
	SP	IF	FA	CO	OT	Failed
BC (one stage)	0	0	0	0	10	10
BC (w/o interact)	8	0	0	0	0	8
BC (w/ food id)	1	0	2	1	1	5
BC (w/ interact)	2	0	2	8	0	10
Rand. Template	2	0	2	0	0	4
Classified Template	10	0	0	0	0	10
MTRNN [1]	0	0	0	0	0	0
SCONE (Ours)	1	0	0	0	0	1

(b) Basic - Sago

Table 1: Basic - Failure Cases

	Extended - Sago					
	SP	IF	FA	CO	OT	Failed
BC (one stage)	0	0	0	0	10	10
BC (w/o interact)	5	1	0	1	0	7
BC (w/ food id)	1	0	2	1	1	5
BC (w/ interact)	0	4	0	3	0	7
Rand. Template	3	4	0	0	0	7
Classified Template	0	0	0	0	0	0
MTRNN [1]	0	5	0	5	0	10
SCONE (Ours)	3	1	0	0	0	4

(a) Extended - Sago

	Peculiar - Penne					
	SP	IF	FA	CO	OT	Failed
BC (one stage)	1	0	8	3	0	9
BC (w/o interact)	1	0	6	1	0	8
BC (w/ food id)	0	0	5	1	0	6
BC (w/ interact)	1	0	0	4	0	5
Rand. Template	4	0	3	0	0	7
Classified Template	2	0	7	0	0	9
MTRNN [1]	2	0	0	0	0	2
SCONE (Ours)	2	0	1	0	0	3

(b) Peculiar - Penne

Table 2: Extended and Peculiar - Failure Cases

Failed Attempts (FA). Failed attempts to scoop are attributed to the same underlying reason as insufficient food. In these cases, the spoon fails to acquire any amount of food, rather than scooping up an insufficient quantity.

Collisions (CO). Relying on vision-based information can lead to incorrect decisions when facing out-of-distribution situations. In such scenarios, collisions between the spoon and the bowl may occur.

Others (OT). In addition to the failure cases mentioned earlier, there are instances where the task cannot be successfully completed or the food ends up being damaged.

4.2 Result Analysis

Table 1: Orange. The BC-based baselines achieved low task success rates due to the high incidence of failed attempts because the weight of the orange used in the testing set is heavier than the foods in the training set. By contrast, the MTRNN and SCONE models exhibited stable performance due to their ability to learn and understand the conditions for successful scooping from the provided demonstrations. The classified template method also achieved higher performance due to its ability to select suitable templates for food items.

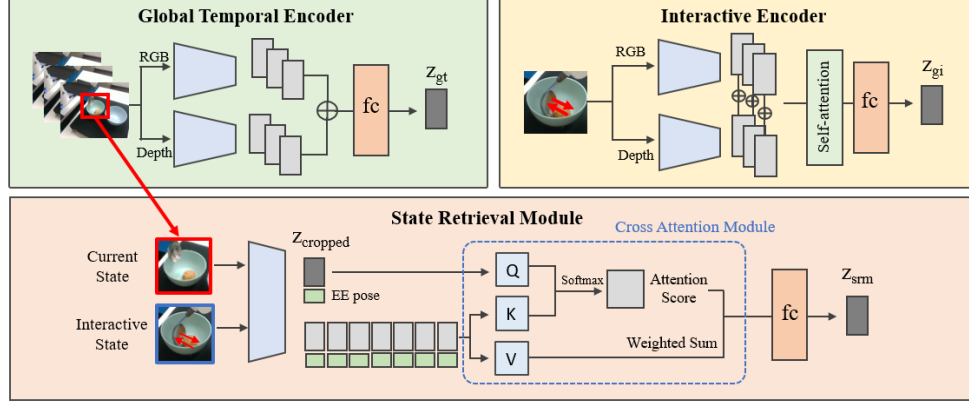


Figure 2: Details of Each Module in SCONE

Table 1: Sago. The BC-based baselines encountered challenges in providing correct predictions, causing several failure cases during the manipulation. Moreover, the template-based method has been limited by its reliance on the predictions from the classifier. In cases where the classifier selected a template intended for food items with large particles, it often resulted in spillage. The MTRNN and the proposed SCONE method demonstrated their ability to overcome this challenge, achieving a higher success rate.

Table 2: Extended Sago. The number of spillages (SP) decreases significantly compared to the basic settings because the amount of sago been reduced, but there was an increase in the occurrence of insufficient food (IF) cases, due to the improper depth insertion of the spoon into the bowl in most of the methods. The Select Template model classified sago as chocolate balls, resulting in 10 successful trials when following the corresponding template. MTRNN demonstrated poor performance under this particular setting. During the scooping stage, it showed a hovering behavior within the bowl and occasionally collided with it. However, SCONE is capable of handling the challenging setting, achieving a success rate of 6 out of 10.

Table 2: Penne. To test the models' generalization abilities, we conducted experiments using peculiar-shaped foods such as penne and fruit candy. Though the color of penne looks familiar to soy beans, the shapes of them are totally different, which led to failed attempts (FA) when employing the soy bean template. While MTRNN was able to handle penne, the jittering trajectory could lead to instability and spillage (SP), which we would like to avoid in real-world evaluation. Our SCONE behaved more stable and maintained a smooth trajectory.

5 Details of Proposed Method

See Figure 2 for detailed information on each encoding module. All the observation inputs are RGB-D images, which are processed through corresponding convolutional layers to extract features.

Global temporal encoder. The global temporal encoder takes as input a sequence of current observations of length N . In our implementation, we set N to 10, which means that the model can access the previous 10 observations within a time window of 1 second. This allows the model to capture the temporal dynamics and dependencies in the input data. Then, the sequences of RGB and depth images are processed separately by their respective encoders, and the output features are concatenated and flattened into a one-dimensional embedding. To further reduce the dimensionality of the features, a fully-connected layer is applied to downsample the features to the dimension of 128.

Interactive encoder. We utilize an interactive encoder to process the sequence of observations captured during the interaction stage. The number of frames K is set to 7. Similar to the encoding process in the global temporal encoder, both RGB and depth images within the sequence are pro-

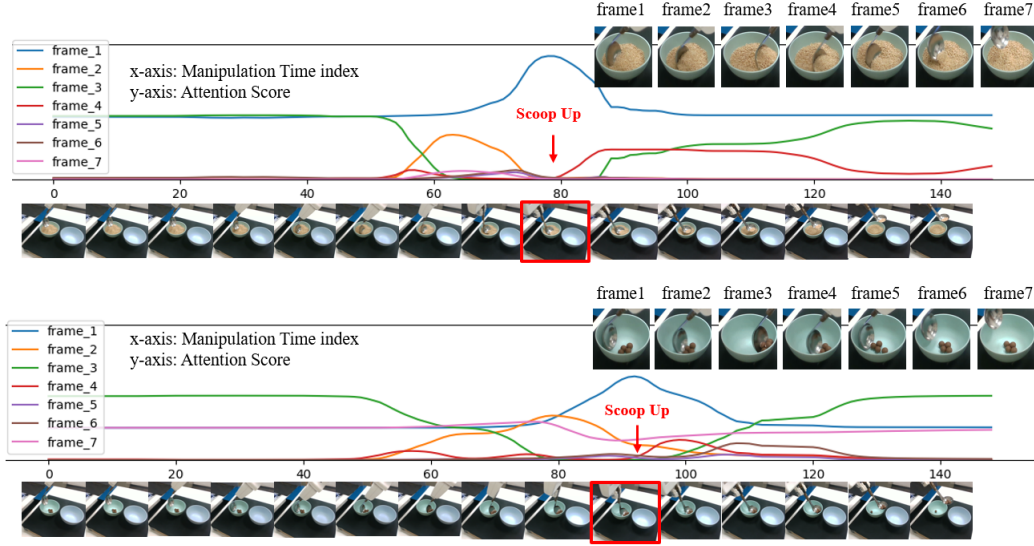


Figure 3: Visualization of Attention Scores. The upper is the manipulation process of brown rice, and the lower is the manipulation process of chocolate ball.

cessed independently through their respective convolutional layers. The output features from these layers are then concatenated to form a sequence of feature maps, capturing the visual information from both modalities. Then, we introduce the multi-head self-attention mechanism, allowing the model to focus on relevant parts of the input sequence. The output features passed through a fully-connected layer, which reduces their dimension to 128.

State retrieval module. The state retrieval module tasks the current local observation and the sequence of interaction as inputs. For the current local observation, which consists of cropped RGB-D images captured at the current time step, we pass them through an encoding module to obtain latent embeddings. These embeddings are then concatenated with the current end-effector pose, resulting in the latent state representation denoted as z_{cropped} . Regarding the sequence of interaction, we apply the same encoding layers to process the observations and obtain corresponding feature embeddings. Similar to the current local observation, we also concatenate the end-effector poses with the encoded features, creating a sequence of states in the interaction. To retrieve the critical state information, we use z_{cropped} as a query and compute its relationship with the sequence of states in the interaction. This is achieved by employing cross-attention mechanisms that output weighted feature embeddings. These embeddings are then flattened and downsampled to a dimension of 32, obtaining the z_{srn} . Furthermore, we also downsample the latent state representation z_{cropped} to a dimension of 32.

6 Visualization of Attention Scores

Figure 3 shows more examples of visualization for attention scores. Based on the results, it is evident that the model can accurately capture the state information without human labeling; this is demonstrated by the consistent patterns observed in the changing attention scores over time at each trial. Overall, our analysis reveals that **frame 1** and **frame 3** exhibit higher attention scores compared to other frames. We attribute this to the presence of important state information related to the spoon’s contact with the food items under specific end-effector poses. These frames can be seen as providing subgoal-like cues, guiding the model in the scooping task. Additionally, we observed that in food items with large particles, **frame 3** tends to have higher attention scores than **frame 1** initially. This is because the state captured in **frame 3** is more similar to the state prior to the scooping action in scenarios involving large particles. On the other hand, we noticed that the highest

177 attention score for **frame 1** typically occurs close to the timing of the actual scooping action. This
178 suggests that **frame 1** serves as a general guide for determining "how" and "where" to scoop up the
179 food items. These findings highlight the model's ability to effectively identify and utilize critical
180 state information from the interacting data to inform its scooping strategy.

181 **References**

- 182 [1] N. Saito, T. Ogata, S. Funabashi, H. Mori, and S. Sugano. How to select and use tools?: Active
183 perception of target objects using multimodal deep learning. *IEEE Robotics and Automation*
184 *Letters*, 6(2):2517–2524, 2021.