

# MoMa-Kitchen: A 100K+ Benchmark for Affordance-Grounded Last-Mile Navigation in Mobile Manipulation

Pingrui Zhang<sup>1,2\*</sup> Xianqiang Gao<sup>2,3\*</sup> Yuhan Wu<sup>3</sup> Kehui Liu<sup>2,4</sup>  
 Dong Wang<sup>2</sup> Zhigang Wang<sup>2</sup> Bin Zhao<sup>2,4</sup> Yan Ding<sup>2†</sup> Xuelong Li<sup>5</sup>

<sup>1</sup>Fudan University <sup>2</sup>Shanghai AI Laboratory <sup>3</sup>University of Science and Technology of China

<sup>4</sup>Northwestern Polytechnical University <sup>5</sup>TeleAI, China Telecom Corp Ltd

{zhangpingrui, dingyan}@pjlab.org.cn

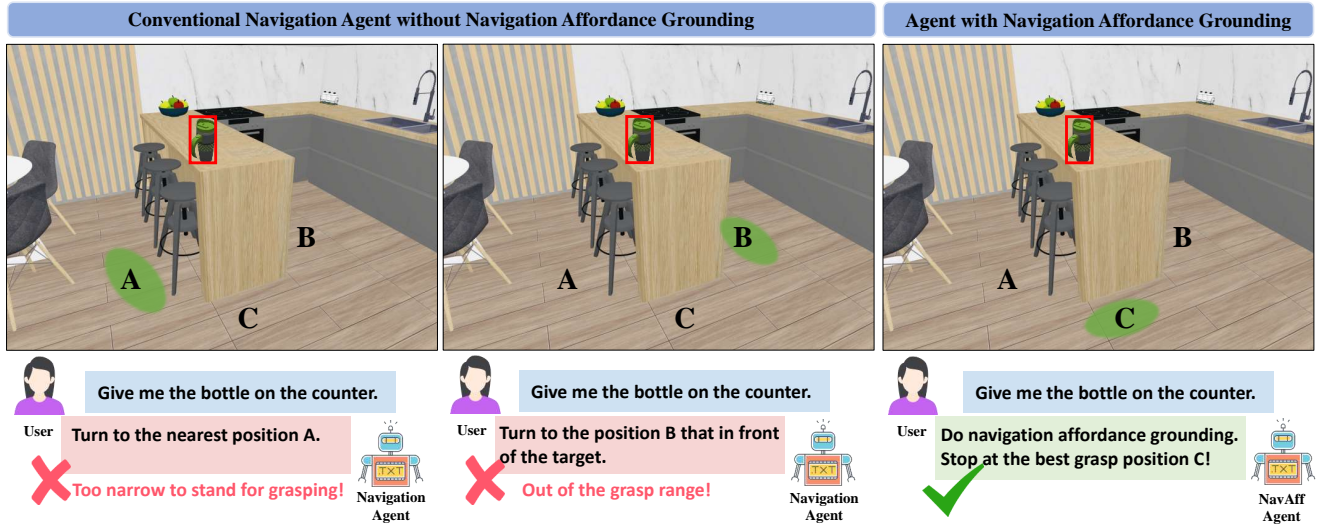


Figure 1. Conventional navigation methods typically prioritize reaching a target location but do not account for constraints affecting manipulation feasibility. *Left:* Position A prioritizes proximity but is obstructed by chairs, preventing stable execution. *Middle:* Position B places the robot in a spacious and stable area for operation but beyond its effective reach. *Right:* Our approach, leveraging navigation affordance grounding, identifies Position C as the optimal stance, ensuring both reachability and task feasibility.

## Abstract

In mobile manipulation, navigation and manipulation are often treated as separate problems, resulting in a significant gap between merely approaching an object and engaging with it effectively. Many navigation approaches primarily define success by proximity to the target, often overlooking the necessity for optimal positioning that facilitates subsequent manipulation. To address this, we introduce MoMa-Kitchen, a benchmark dataset comprising over 100k samples that provide training data for models to learn optimal final navigation positions for seamless transition to manipulation. Our dataset includes affordance-grounded floor labels collected from diverse kitchen environments, in which robotic mobile manipulators of different models attempt to

grasp target objects amidst clutter. Using a fully automated pipeline, we simulate diverse real-world scenarios and generate affordance labels for optimal manipulation positions. Visual data are collected from RGB-D inputs captured by a first-person view camera mounted on the robotic arm, ensuring consistency in viewpoint during data collection. We also develop a lightweight baseline model, NavAff, for navigation affordance grounding that demonstrates promising performance on the MoMa-Kitchen benchmark. Our approach enables models to learn affordance-based final positioning that accommodates different arm types and platform heights, thereby paving the way for more robust and generalizable integration of navigation and manipulation in embodied AI. Project page: <https://momakitchen.github.io/>.

# 1. Introduction

Most existing navigation algorithms define success *in terms of* reaching a location near the target [47]. However, in household environments, navigation is an intermediate step preceding task-specific manipulation. As a result, such navigation strategies are inadequate for mobile manipulation tasks that requires precise end-effector positioning. In practice, navigation and manipulation are tightly integrated. For instance, when tasked with retrieving an object from a kitchen counter, a robot typically navigates toward the target using conventional policies before attempting manipulation. While existing navigation algorithms reliably guide robots across rooms, they often fail in proximity to the target. The robot may stop beyond the manipulator’s reachable workspace or be obstructed by spatial constraints, rendering manipulation infeasible (as shown in Fig. 1). Additionally, obstacles such as furniture, bins, or containers in cluttered environments are seldom accounted for in existing navigation policies, further complicating the selection of feasible grasping positions. As a result, reliance solely on these navigation algorithms, without incorporating the demands of manipulation, limits their efficacy in addressing complex tasks in household settings.

This limitation highlights the disconnect between object localization (*navigation*) and physical interaction (*manipulation*) in mobile robotics. In particular, optimizing final positioning in the “last mile” remains a fundamental challenge, yet existing datasets and benchmarks provide limited supervision for this aspect. While recent efforts have employed large language models (LLMs) to assist in selecting optimal navigation positions, these approaches fall short when transitioning to the manipulation phase [44, 86, 88]. Training-free LLMs struggle to accurately predict the requirements of robotic arm interactions. Additionally, they cannot dynamically adjust positioning strategies based on different robotic arm models or base morphology [72, 87].

To bridge this gap, we introduce **MoMa-Kitchen**, a large-scale benchmark with over 100k samples designed to train models for affordance-grounded final positioning in mobile manipulation tasks. Our dataset comprises 127,343 episodes spanning 569 diverse kitchen scenes, where each episode involves predicting floor affordances that enable a robot to approach a target while avoiding collisions with obstacles. We collect large-scale floor affordance data by creating kitchen scenes with various layouts—in which target objects are either randomly placed or selected from common appliances (e.g., refrigerators and cabinets) and obstacles are strategically positioned to generate distinct scenarios. In the simulator, various mobile manipulators attempt to grasp target objects from multiple positions, and success rates are recorded to obtain ground truth affordance labels for each floor position. The collected visual data, consisting of RGB-D and point cloud inputs from a first-

person view camera mounted on the robotic arm, along with robot-related information, are subsequently used to train our lightweight baseline model, **NavAff**, which identifies the optimal floor affordance region for subsequent manipulation. To enhance generalization, we employ various robotic arms (e.g., Flexiv and Franka) and varied mobile bases during grasping, enabling the model to learn that floor affordance predictions depend on arm height and operational range. This approach is intended to develop a model that can be used on heterogeneous devices [60]. In summary, our contributions are as follows:

- We propose MoMa-Kitchen, the first large-scale dataset with over 100k samples that bridges the gap between navigation and manipulation in mobile manipulation tasks by enabling models to optimize final positioning near target objects.
- We develop a fully automated data collection pipeline – including scene generation, affordance labeling, and object placement – to simulate diverse real-world scenarios and enhance model generalizability.
- We design a lightweight baseline model, NavAff that employs RGB-D and point cloud inputs for navigation affordance grounding, achieving promising results on the MoMa-Kitchen benchmark.

## 2. Dataset Generation

**Task Definition.** MoMa-Kitchen focuses on determining feasible final navigation positions that enable successful manipulation in cluttered environments. Given RGB-D inputs from a first-person camera and robot-specific parameters (e.g., arm reach, base height), the goal is to produce an affordance map over the floor. This map indicates where the robot can position itself to reliably manipulate the target while accounting for obstacles, bridging navigation and manipulation within a unified pipeline.

MoMa-Kitchen includes diverse kitchen environments with multiple types of robotic mobile manipulators, first-person view visual data, and navigation-specific affordance ground truth. As illustrated in Fig. 2, we first construct large-scale kitchen environments that contain different targets, obstacles, and furniture. Mobile manipulators are then placed in these environments to obtain navigation affordance labels through robotic manipulation. For each labeled scene, we collect first-person view visual data from multiple randomly sampled robot viewpoints, including RGB-D images of the scene. The position of each viewpoint is also recorded to generate and transform both the global and floor point clouds. Additionally, floor-level affordance ground truth is collected near the target.

### 2.1. Scene Setup

To construct diverse kitchen environments for affordance annotation, we employ BestMan [74], a PyBullet-based

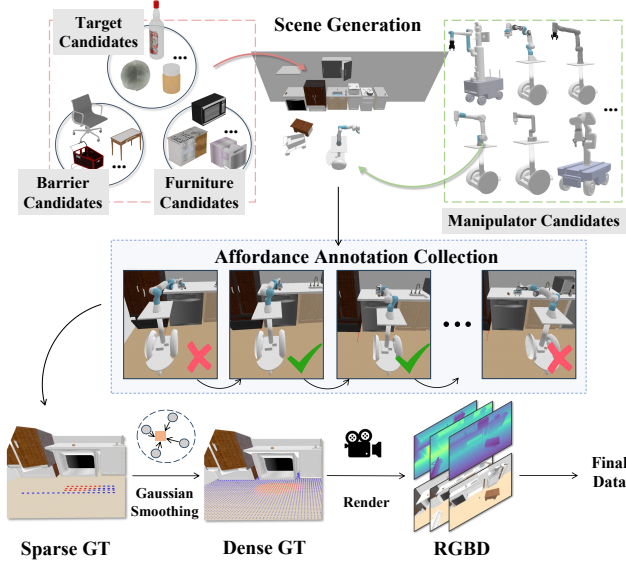


Figure 2. **Overview of scene setup and data generation pipeline.** Each scene features unique base furniture and layout, with randomly placed obstacles surrounding the target object to enhance scene complexity. Discrete navigation affordance values are collected by moving the mobile manipulator and interacting with the target objects in the scene. View transformation and Gaussian interpolation are then applied to generate a dense affordance map, along with corresponding RGBD data.

simulation platform that integrates assets from PartNet-Mobility [8, 38, 71]. The generation process begins with a rectangular kitchen layout, where common kitchen furniture and appliances (e.g., sinks, cabinets, dishwashers, and fridges) are procedurally arranged along one wall. To ensure scene diversity, we randomize both the placement order of object categories and the specific instance selection within each category. We further augment each scene with both rigid and articulated objects as manipulation targets, and introduce additional obstacles around them to increase scene complexity. Once the scene assets are ready, we randomly select a manipulator from various types and place it into the scene to collect navigation affordance labels.

## 2.2. Visual Data Collection

To comprehensively capture the target and its surroundings from the robot’s perspective, we sample ten distinct camera poses  $[\mathbf{R}_c | \mathbf{T}_c] \in \mathbb{R}^{4 \times 4}$  around each target. These viewpoints are strategically selected to cover the target object, its surrounding environment, and the floor. For viewpoint selection, we alternate the positions to the left and right of the target object, gradually increasing the distance from it. After placing the robot at the selected position, we check whether the target object is within the observation range. If not, we randomly generate a new position. This process continues until the target object is within the line of sight,

at which point we stop and record the viewpoint. For each camera pose, we collect RGB images  $\mathbf{I}$ , depth maps  $\mathbf{D}$ , and floor point clouds  $\mathbf{P}_{\text{floor}} \in \mathbb{R}^{n \times 3}$ , where  $n$  denotes the number of points in the point cloud.

## 2.3. Affordance Labeling

We gather mobile manipulators with various robotic arms to collect floor-level navigation affordance ground truth data for target objects. For each target object, we define the affordance sampling area  $\mathcal{A}$  as a semicircular region on the floor, centered at the target’s position with a radius equal to the maximum reach of the robotic arm. At each sampling position  $\mathbf{p} \in \mathcal{A}$ , we place the robot base and initialize the manipulator configuration, then align the end-effector orientation  $\mathbf{R} \in \text{SO}(3)$  with the surface normal of either the target object (for rigid objects) or its functional link (for articulated objects). Here,  $\text{SO}(3)$  denotes the Special Orthogonal Group, representing the set of all rotation matrices that describe rotations in three-dimensional space.

For each configuration  $(\mathbf{p}, \mathbf{R})$ , we attempt a manipulation and record a binary affordance outcome  $v_p \in \{0, 1\}$  at position  $\mathbf{p}$ . Success is determined by the outcome of the manipulation attempt: for robots equipped with two-finger parallel grippers (e.g., Panda, Flexiv, and xArm6), success is defined by a successful grasp of the target object; for the robot employing a suction-based end-effector (e.g., UR5e), success is determined by a valid suction-and-move action. To associate these affordance values with the previously collected floor point cloud  $\mathbf{P}_{\text{floor}}$ , we first transform  $v_p$  from the world coordinate system to the robot base frame. This transformation yields  $v_{p_{\text{base}}}$  through:

$$v_{p_{\text{base}}} = \mathbf{T}_{bc} \mathbf{T}_{cw} v_p, \quad (1)$$

where  $\mathbf{T}_{bc}$  and  $\mathbf{T}_{cw} \in \mathbb{R}^{4 \times 4}$  are the transformations from the camera to the base and from the world to the camera, respectively. We then match each transformed affordance value  $v_{p_{\text{base}}}$  to its nearest neighbor in  $\mathbf{P}_{\text{floor}}$ . This association process can be formally expressed as:

$$\mathbf{V}_{\text{aff}}(p_j) = \begin{cases} v_{p_{\text{base}}} & \text{if } \min_{p \in \mathcal{A}} \|p - p_j\| < \theta, \forall p_j \in \mathbf{P}_{\text{floor}}, \\ 0 & \text{otherwise,} \end{cases} \quad (2)$$

where  $\theta$  is a distance threshold. This yields sparse affordance values  $\mathbf{V}_{\text{aff}} \in \mathbb{R}^{n \times 1}$ . To produce dense and continuous floor-level affordance maps, we employ Gaussian interpolation with  $k$ -nearest neighbors. This process produces interpolated affordance values  $\hat{\mathbf{V}}_{\text{aff}} = \{\hat{v}_i \in [0, 1] \mid i = 1, \dots, n\}$  through the following formulation:

$$\hat{v}_i = \frac{\sum_{j=1}^k w_{ij} v_j}{\sum_{j=1}^k w_{ij}} \quad (3)$$

where the Gaussian weights  $w_{ij}$  are computed as follows.

$$w_{ij} = \exp\left(-\frac{\|p_i - p_j\|^2}{2\sigma^2}\right), \quad \forall i, j \quad (4)$$

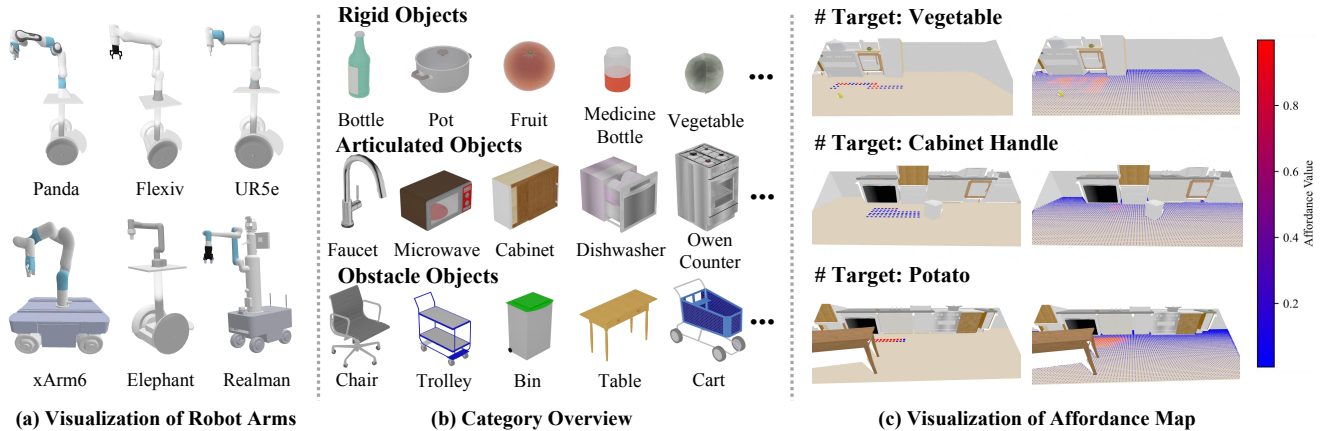


Figure 3. **(a) Robot arms used in MoMa-Kitchen.** The end-effectors of the Panda, Flexiv, Elephant, Realman and xArm6 robots are grippers, while the end-effector of the UR5e is a suction cup. **(b) Object categories utilized in MoMa-Kitchen.** Each category consists of multiple object instances. Rigid and articulated objects serve as manipulation targets, while obstacle objects are strategically placed around the target to enhance scene complexity. **(c) Examples of affordance maps in MoMa-Kitchen.** Discrete affordance values are first collected (left) by moving the mobile manipulator and allowing it to interact with the target. Gaussian interpolation is then applied to obtain a smooth affordance map (right).

Here,  $\|p_i - p_j\|$  is the Euclidean distance between the interpolation target  $p_i$  and the sparse point  $p_j$ ,  $\sigma$  controls the width of the Gaussian kernel, and  $s_j$  is the affordance value at  $p_j$ . This weighting scheme ensures that points closer to  $p_i$  have a larger influence, resulting in a smooth, continuous affordance map that captures the spatial distribution of potential manipulation outcomes.

## 2.4. Generated Dataset Statistics

As summarized in Table 1, MoMa-Kitchen spans 569 kitchen scenes, each with procedural variations in objects and obstacles. These configurations yield over 127k episodes, each capturing an RGB-D view and the resulting affordance map. Six robot arms (Flexiv, Panda, UR5e, xArm6, Realman, Elephant) are deployed to ensure diverse reachability constraints and end-effector types. (see Fig. 3(a)) The dataset covers 137 kitchen-relevant assets (65 rigid, 48 articulated, and 24 obstacle types), providing a broad spectrum of layouts and manipulation targets. (see Fig. 3(b)) (Details in Supp.1) During data collection, each configuration is labeled via discrete success/failure outcomes for different robot placements. These discrete labels are then interpolated to produce dense affordance maps (see Fig. 3(c)). Collectively, MoMa-Kitchen provides large-scale supervision for end-to-end training of navigation-to-manipulation models that generalize across varying hardware and scene complexity.

## 3. Baseline Model

In this section, we introduce our baseline model, NavAff, which is designed for optimal manipulation positioning and

Table 1. **Dataset Split Statistics.**

Split	Scenes	Configurations	Episodes
Train	456	11,408	102,687
Test (Unseen Scenes)	113	2,747	24,656

obstacle interaction in complex mobile manipulation tasks. As shown in Fig. 4, Our baseline method consists of two main components: Visual Alignment Module (VAM) and Navigation Affordance Grounding Module (NAG). This two-part structure is designed to process visual and spatial information sequentially, facilitating effective affordance prediction for mobile manipulation tasks.

### 3.1. Visual Alignment

VAM extracts and projects features of target object  $T$  and obstacles  $\{O_1, O_2, \dots, O_n\}$  onto the global point cloud  $\mathbf{P}_{global} \in \mathbb{R}^{3 \times N_{global}}$ , to ensure accurate 3D representation of them. As illustrated in Fig. 4(a), VAM processes an image  $\mathbf{I} \in \mathbb{R}^{3 \times H \times W}$  captured by the robot’s camera along with its corresponding depth  $\mathbf{D} \in \mathbb{R}^{H \times W}$ , extracting target and obstacle masks using BestMan[74], where  $H, W$  denotes the height and width of the camera observation. The global point cloud  $\mathbf{P}_{global}$  is then generated via inverse projection. To enhance feature representation, we leverage PointNet++ [41], which supports additional feature channels in the point cloud input. Pixels from the target mask are projected into the 3D space to form a *target channel*, with values at target locations set to 1, while obstacle mask pixels create an *obstacle channel*, assigned related values of -1 [22]. This process enriches the global

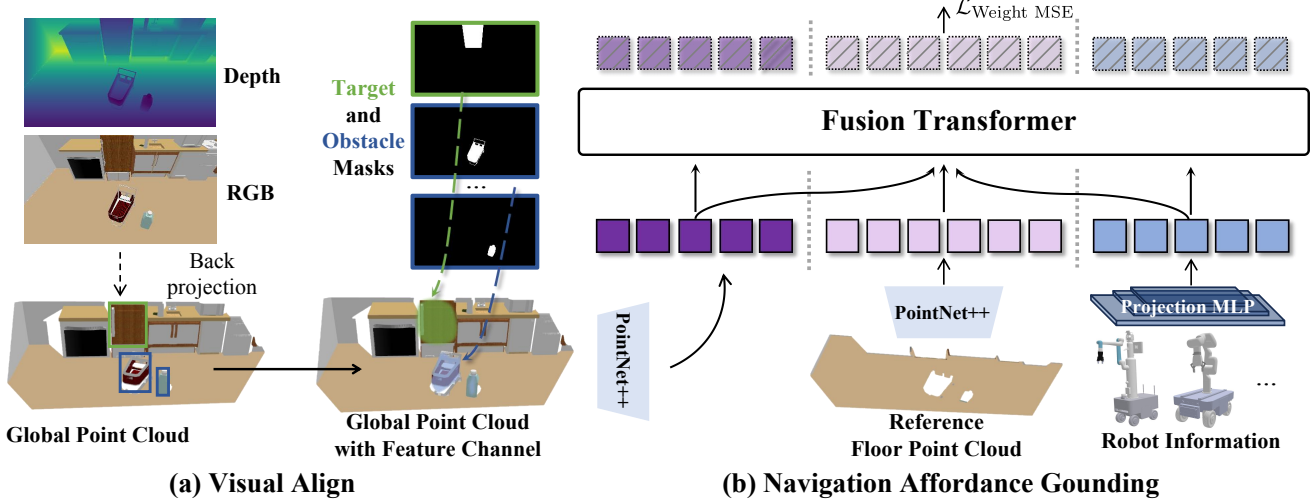


Figure 4. **NavAff Baseline.** (a) **Visual Alignment Module:** Projects object masks to align 2D visual features with 3D spatial representations. (b) **Navigation Affordance Grounding Module:** Fuses global point cloud, floor point cloud, and robot-specific features to predict navigation affordance maps.

point cloud by embedding 2D visual cues (e.g., target and obstacle features) into the 3D spatial structure. The resulting feature-enhanced point cloud  $\bar{\mathbf{P}}_{global}$  serves as input for the next stage, aligning visual perception with spatial information for improved downstream processing.

### 3.2. Navigation Affordance Grounding

NAG uses the robot’s relevant information and the visually aligned global point cloud to interact with the reference floor, generating a robot-specific affordance prediction for optimal positioning. As illustrated in Fig. 4(b), NAG takes as input the feature-enhanced global point cloud  $\bar{\mathbf{P}}_{global}$  from VAM, the floor point cloud  $\mathbf{P}_{floor}$ , and the robot-related information  $\mathbf{R}_I$ . To extract features from the point cloud data, PointNet++ processes  $\bar{\mathbf{P}}_{global}$  and  $\mathbf{P}_{floor}$ , yielding feature representations  $\mathbf{F}_{pcg}$  and  $\mathbf{F}_{pcf}$ , respectively. These features are then tokenized using a multi-layer perceptron (MLP) in preparation for multi-head cross-attention [56]:

$$\bar{\mathbf{F}}_{pcg} = \text{Tokenizer}(\mathbf{F}_{pcg}), \quad (5)$$

$$\bar{\mathbf{F}}_{pcf} = \text{Tokenizer}(\mathbf{F}_{pcf}), \quad (6)$$

where  $\bar{\mathbf{F}}_{pcg} \in \mathbb{R}^{n_{global} \times dim}$  and  $\bar{\mathbf{F}}_{pcf} \in \mathbb{R}^{n_{floor} \times dim}$ .

For the mobile manipulator,  $\mathbf{R}_I$  encodes the base platform height and the operational radius of the robotic arm. This information is also processed by an MLP to produce the robot-specific tokens  $\mathbf{F}_r$ , which are then concatenated with  $\bar{\mathbf{F}}_{pcg}$ , forming a combined key and value for cross-attention fusion with  $\bar{\mathbf{F}}_{pcf}$ . The process is defined as follows:

$$\mathbf{F}_r = \text{Tokenizer}(\text{MLP}(\mathbf{R}_I)) \quad (7)$$

$$\mathbf{F}_{out} = \text{MCA}(\bar{\mathbf{F}}_{pcf} \mathbf{W}_q, [\bar{\mathbf{F}}_{pcg}, \mathbf{F}_r] \mathbf{W}_k, [\bar{\mathbf{F}}_{pcg}, \mathbf{F}_r] \mathbf{W}_v) \quad (8)$$

where  $\mathbf{F}_r \in \mathbb{R}^{n_{robot} \times dim}$ ,  $\mathbf{F}_{out} \in \mathbb{R}^{n_{floor} \times dim}$ . Finally,  $\mathbf{F}_{out}$  is passed through a decoder  $f_d$  to predict the navigation affordance grounding  $\mathbf{P}_{out}$ .

### 3.3. Loss Function

Following standard practices in affordance grounding, we utilize the mean squared error (MSE) loss as the primary objective function to align the model’s predictions  $\mathbf{P}_{out}$  with the ground truth affordance labels  $\mathbf{P}_{label}$ . In our experiments, we observed that zero-valued elements constitute a large proportion of the ground truth floor affordance labels, leading to an imbalance between zero-valued elements and those with non-zero values. To address this issue, we apply a weighted mask to perform a weighted average on the MSE loss (Weighted MSE). The Weighted MSE loss is computed by adjusting the weight for zero-valued elements in the ground truth. Specifically, for elements in the ground truth that are zero, we apply a weight of  $\lambda$  with a probability of 0.5, while other elements are assigned a weight of 1. This ensures the model places adequate emphasis on non-zero affordance areas. The formula for the Weighted MSE loss is as follows:

$$\mathcal{L}_{\text{Weight MSE}} = \frac{1}{N} \sum_{i=1}^N \mathbf{W}_i \cdot (\mathbf{P}_{out,i} - \mathbf{P}_{label,i})^2 \quad (9)$$

where

$$\mathbf{W}_i = \begin{cases} \lambda & \text{if } \mathbf{P}_{label,i} = 0 \text{ and with prob. } 0.5 \\ 1 & \text{otherwise.} \end{cases} \quad (10)$$

Here,  $N$  represents the total number of elements,  $\mathbf{W}_i$  is the weight applied to each element based on the corresponding floor ground truth, and  $\lambda \in (0, 1)$  is a hyperparameter that assigns a reduced weight to samples with a zero ground truth value, applied with a probability of 0.5.

## 4. Experiments

In this section, we benchmark NavAff on the MoMa-Kitchen dataset using the data split described in Tab. 1, with affordance annotations from Sec. 2.3 and visual data from Sec. 2.2. We introduce baseline results for NavAff, conduct a comprehensive performance analysis across both simulated and real-world environments, and identify emerging challenges in complex mobile manipulation scenarios that are uniquely highlighted through our benchmark.

Table 2. **Main results.** Quantitative evaluation of navigation affordance grounding performance of NavAff.

Method	RMSE ↓	logMSE ↓	PCC ↑	SIM ↑
PointNet++ [41]	0.164	0.0142	0.565	0.589
VoteNet [42]	0.167	0.0143	0.543	0.570
H3DNet [83]	0.174	0.0156	0.503	0.522
NavAff	<b>0.147</b>	<b>0.0115</b>	<b>0.680</b>	<b>0.696</b>

### 4.1. Experimental Settings

Given that this is a newly proposed task, no existing methods provide a direct basis for comparison. To establish an evaluation, we adapt three well-established models from point cloud learning: PointNet++[41], VoteNet [42], and H3DNet [83]. Specifically, PointNet++ serves as a foundational backbone model widely used in point cloud processing tasks; VoteNet and H3DNet, both originally designed for 3D object detection, are suited for adaptation to navigation affordance grounding on our benchmark. By leveraging these models, we aim to establish strong baselines and gain insights into how existing point cloud techniques perform when applied to this new challenge.

For evaluation, we employ standard and diversity metrics tailored for the MoMa-Kitchen benchmark. Each metric provides a distinct perspective on navigation affordance grounding performance: Root Mean Squared Error (RMSE), facilitating understanding the magnitude of prediction error; Logarithmic Mean Squared Error (logMSE) focuses more on relative differences rather than absolute differences; Pearson Correlation Coefficient (PCC) helping to gauge the model’s ability to maintain consistent patterns with the ground truth data across various affordance regions; Cosine Similarity (SIM) compares the structural or shape similarity between predicted and ground truth.

All the above experiments are trained on a single NVIDIA A100 GPU with a batch size of 64, using the

Table 3. **Manipulation Success Rate (MSR).**

MSR	Random	H3DNet	VoteNet	PointNet++	NavAff
Top1	0.080	0.54	0.56	0.60	<b>0.72</b>
Top5	0.046	0.47	0.53	0.58	<b>0.66</b>

Adam optimizer with a learning rate of  $8e-4$ . Further experimental details are available in the Appendix.

### 4.2. Quantitative Analysis

#### 4.2.1. Main Results

Tab. 2 reports the metrics of the proposed method NavAff compared with other methods on the MoMa-Kitchen benchmark. The results highlight that NavAff achieves superior performance in navigation affordance grounding, outperforming transferred methods across all evaluated metrics. Specifically, NavAff achieves an RMSE of 0.147 and a logMSE of 0.0115, demonstrating a marked improvement in prediction accuracy compared to the second-best method, PointNet++. In terms of correlation and similarity measures, NavAff achieves the highest scores with a PCC of 0.680 and SIM of 0.696, highlighting its robustness in capturing affordance patterns even in cluttered environments. Notably, VoteNet and H3DNet fall short of NavAff’s accuracy and consistency across the various metrics, showcasing the effectiveness of our approach in the proposed MoMa-Kitchen benchmark.

For further analysis, while VoteNet and H3DNet are effective in point cloud detection and proposal classification tasks, they are less suited to the fine-grained requirements of navigation affordance grounding, resulting in slightly lower performance. All methods are trained for the same number of epochs, but H3DNet converges more slowly due to its larger parameter count. Consequently, within the same training duration, H3DNet demonstrates the weakest performance. These findings suggest that a lightweight model capable of fine-grained feature prediction may achieve superior results on our MoMa-Kitchen benchmark.

#### 4.2.2. Manipulation Success Rate

To more intuitively demonstrate the direct improvement of our method and benchmark its performance in mobile manipulation, we introduce the manipulation success rate (MSR). The results are shown in Tab. 3. The Top1 MSR refers to moving the robot to the location with the highest affordance score and recording the MSR in the test scenes. The Top5 MSR refers to recording the average MSR across the top 5 locations.

As quantitatively shown in Tab. 3 using the MSR metric, our navigation affordance prediction paradigm shows significant practical advantages. A comparison between the “Random” baseline and other methods reveals that approaches trained on MoMa-Kitchen benchmark signif-

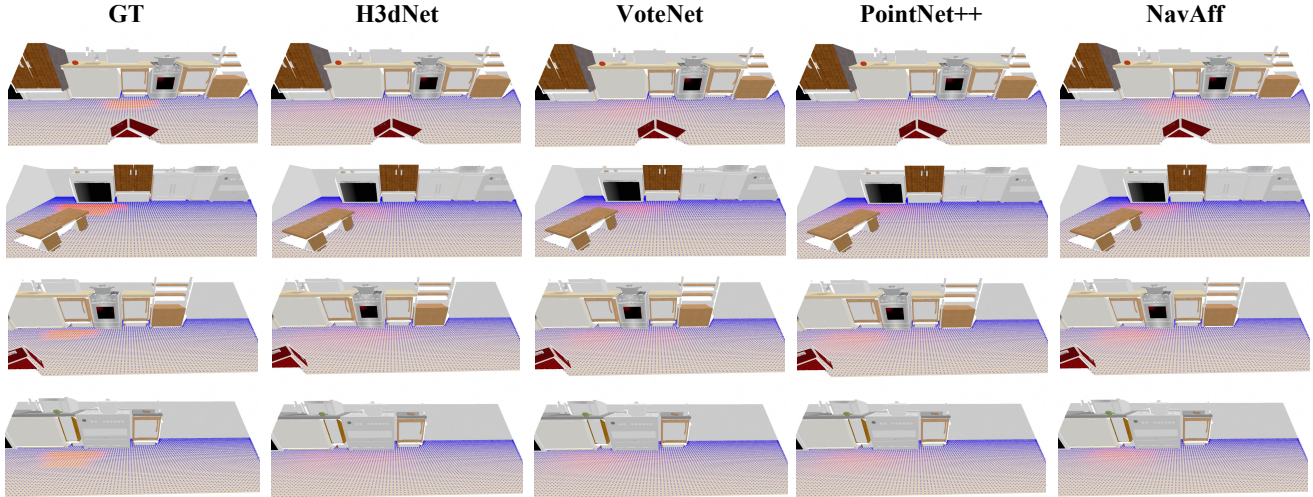


Figure 5. **Qualitative comparison of navigation affordance between all methods and ground truth.** Blue to red regions indicate affordance values ranging from 0 to 1, while void areas represent obstacle-occupied spaces.

icantly outperform those relying on randomly sampled points within the robot’s operational radius, thereby validating the effectiveness of our benchmark. Among all the trained methods, NavAff sets a new state-of-the-art performance, achieving an impressive  $\text{Top1 MSR}$  of 72% and  $\text{Top5 MSR}$  of 66%. Notably, when expanding the candidate point selection from Top1 to Top5, the model demonstrates strong generalization with minimal performance degradation, suggesting its robust adaptability to varying operational conditions.

Additionally, by comparing the MSR values of all models, we observe a strong positive correlation between MSR and the accuracy of navigation affordance prediction (see Tab. 2 and Tab. 3), indicating that accurate navigation affordance prediction indeed contributes to higher manipulation success rates in mobile manipulation tasks.

### 4.3. Qualitative Results

We show the qualitative results of NavAff and compared methods for navigation affordance grounding in Fig. 5. As illustrated, our model successfully predicts floor-level affordance maps that closely align with the ground truth patterns. The visualization clearly reveals distinct void regions in both predicted and ground truth maps, which correspond to areas occupied by obstacles, effectively demonstrating our model’s ability to recognize spatial constraints. Compared to other methods, NavAff exhibits a more accurate representation of these void regions, indicating its superior capability in handling complex spatial relationships and navigating cluttered environments.

### 4.4. Ablation Study

To understand the impact of each module in our proposed NavAff model, we conduct an ablation study by progres-

sively removing components and evaluating performance on the MoMa-Kitchen benchmark. Tab. 4 shows the results of this study, where we assess the model variations by excluding robot information, the Visual Alignment Module (VAM), and the global point cloud. Below, we analyze the role and impact of each component in detail.

- **w/o robot information.** Removing robot-specific parameters (*e.g.*, base height and arm reach) leads to moderate performance degradation, particularly in RMSE and logMSE. The limited impact may be attributed to the simplified representation of robot information in NavAff, which only considers two parameters. For instance, only xArm6 has a unique base height, while other platforms share identical base configurations. Future work could explore incorporating additional robot-specific parameters to enhance performance.
- **w/o VAM.** Excluding the Visual Alignment Module significantly reduces performance, particularly in PCC and SIM metrics, underscoring the importance of integrating 2D visual cues for accurate navigation affordance prediction. This demonstrates that aligning visual data with spatial information is essential for accurate and robust navigation affordance grounding.
- **w/o global point cloud.** Without the global point cloud, the model’s performance drops across all metrics, with the largest degradation in PCC. This indicates that global spatial context is crucial for capturing the layout of objects and obstacles in the scene.

### 4.5. Real World Application

As shown in Fig. 6, we validate our method by conducting experiments in a real-world setting. A D435i camera is mounted on the head of a mobile manipulator to capture data. Using Grounded-Sam [46], we obtain open-

Table 4. Ablation study results of NavAff.

	RMSE ↓	logMSE ↓	PCC ↑	SIM ↑	top1MSR ↑
NavAff	0.147	0.0115	0.680	0.696	0.72
w/o robot information	0.148	0.0115	0.670	0.688	0.70
w/o VAM	0.165	0.0140	0.562	0.589	0.63
w/o global point cloud	0.168	0.0144	0.534	0.568	0.58

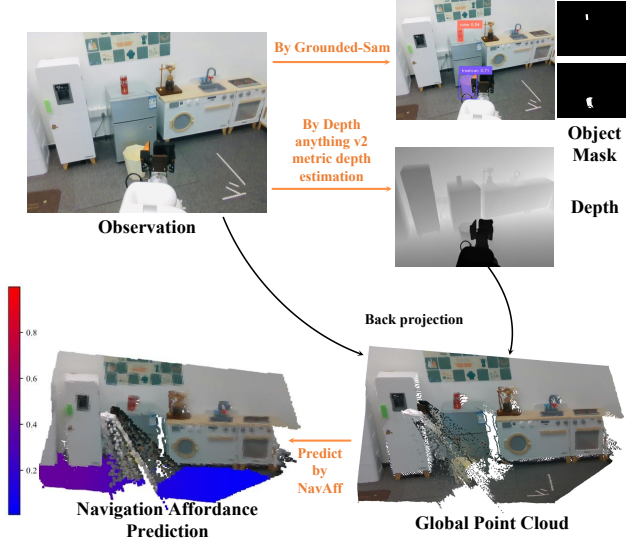


Figure 6. Pipeline of real-world application with simulator trained NavAff. The D435i camera mounted on the mobile manipulator captures image, and our model predicts optimal navigation affordances on the floor based on this data.

vocabulary object masks, while Depth Anything v2 [76] provided the depth images. We generate the global point cloud through back projection and segmented the object’s local point cloud based on the mask. Then, we apply our model—trained on simulation data—to predict navigation affordances on the floor.

The prediction results, shown in Fig. 6, demonstrate that our model performs well in real-world kitchen scenarios, validating the strong generalization capability of our benchmark. This success highlights the robustness of our method, which was specifically designed to minimize the visual gap between simulation and reality. By addressing this challenge, our model is able to effectively adapt to real-world settings, showing that our approach is not only effective in controlled environments but also highly reliable in dynamic, real-world situations. Details of mobile manipulation demo are available in the supplementary material.

## 5. Related Work

**Mobile Manipulation.** Such tasks have been extensively studied by researchers in both simulated and real-world settings [16, 17, 37, 49, 67, 68, 78]. More recently, approaches that integrate visual and linguistic information have emerged as promising methods for achieving unified reasoning in mobile manipulation [33, 44, 58, 77, 85]. How-

ever, most existing methods still lack the versatility required to seamlessly combine coarse and fine motions for both navigation and manipulation. Prior mobile manipulation approaches often assume obstacle-free or easily navigable environments. While interactive navigation techniques [70] attempt to tackle scenarios in which obstacles must be moved or manipulated—such as shifting boxes or pressing buttons—they typically rely solely on geometric reasoning [30, 50, 55, 61, 63, 80]. In contrast, our proposed benchmark is specifically designed to address these limitations by integrating both semantic and geometric reasoning for optimal positioning and manipulation. This benchmark provides a robust and cohesive foundation for advancing mobile manipulation tasks in complex environments.

**Embodied Dataset.** In recent years, large-scale embodied datasets have become essential for advancing both navigation and manipulation tasks in robotics. On navigation, several datasets have been proposed to tackle challenges in object-goal navigation [7, 25, 45, 69], rearrangement [4, 36, 65], vision-language navigation [3, 27, 28, 43], and question answering [10, 11, 66, 79, 82]. In contrast, manipulation-focused works leverage large-scale proprioceptive and visual data (e.g., RGB images or 3D point clouds) to simulate real-world dexterous hand operations [31, 34, 57, 62, 73], while other datasets provide object affordance labels to facilitate interaction learning [5, 9, 15, 18, 23, 52, 75, 81]. In parallel to simulation-based approaches, several real-world datasets have emerged to address the gap between simulated and physical environments. Some efforts like RoboTurk [35], MIME [48], RoboNet [12] provide valuable demonstrations of physical object interactions, and FastUMI [84] propose a scalable hardware-independent robotic manipulation data collection system, while datasets such as KITTI [19], nuScenes [6], and Waymo Open [51] have significantly advanced autonomous navigation. More recently, research has increasingly focused on integrating manipulation capabilities with navigation for everyday tasks. Recent frameworks like BEHAVIOR Robot Suite [24] enable whole-body manipulation with bimanual coordination in household environments, while AgiBot World Colosseo [1] offers over one million trajectories across 217 tasks. Despite these advances, real-world datasets still face challenges in human supervision and scalability [29, 39]. Moreover, some approaches use vision-language models or procedural methods to autonomously generate scalable language annotations in simulated environments [2, 13, 14, 20, 21, 32, 53, 54, 59, 64]. Despite these advances, navigation datasets often offer rich spatial information yet lack guidance for optimal positioning during subsequent manipulation, whereas manipulation datasets—despite providing valuable interaction data—do not fully capture the complexities of achieving optimal grasping positions via navigation.



## 6. Conclusion

Our work addresses the “last mile” challenge in mobile manipulation by introducing MoMa-Kitchen, the first large-scale dataset (127k+ episodes across 569 kitchen scenes) featuring comprehensive, high-quality ground truth affordance maps to guide optimal positioning for manipulation tasks. We resolve the navigation-manipulation disconnect through automated cross-platform data collection and a unified framework compatible with diverse robotic systems, ensuring generalizability. Experimental results demonstrate that our baseline model NavAff achieves robust affordance prediction, validating the approach’s efficacy. This dataset and methodology not only advance integrated navigation-manipulation systems for real-world deployment but also provide essential infrastructure for scalable and adaptive robotic learning in dynamic, cluttered environments.

## References

- [1] AgiBot-World-Contributors, Qingwen Bu, Jisong Cai, Li Chen, Xiuqi Cui, Yan Ding, Siyuan Feng, Shenyuan Gao, Xindong He, Xu Huang, Shu Jiang, Yuxin Jiang, Cheng Jing, Hongyang Li, Jialun Li, Chiming Liu, Yi Liu, Yuxiang Lu, Jianlan Luo, Ping Luo, Yao Mu, Yuehan Niu, Yixuan Pan, Jiangmiao Pang, Yu Qiao, Guanghui Ren, Cheng-Xing Ruan, Jiaqi Shan, Yongjian Shen, Chengshi Shi, Mi Shi, Modi Shi, Chonghao Sima, Jianheng Song, Huijie Wang, Wenhao Wang, Dafeng Wei, Chengen Xie, Guo-Liang Xu, Junchi Yan, Cunbiao Yang, Lei Yang, Shukai Yang, Maoqing Yao, Jia Zeng, Chi Zhang, Qingli Zhang, Bin Zhao, Chengyu Zhao, Jiaqi Zhao, and Jianchao Zhu. Agibot world colosseum: A large-scale manipulation platform for scalable and intelligent embodied systems. 2025. 8
- [2] Michael Ahn, Debidatta Dwibedi, Chelsea Finn, Montse Gonzalez Arenas, Keerthana Gopalakrishnan, Karol Hausman, Brian Ichter, Alex Irpan, Nikhil Joshi, Ryan Julian, et al. Autort: Embodied foundation models for large scale orchestration of robotic agents. *arXiv preprint arXiv:2401.12963*, 2024. 8
- [3] Peter Anderson, Qi Wu, Damien Teney, Jake Bruce, Mark Johnson, Niko Sünderhauf, Ian Reid, Stephen Gould, and Anton Van Den Hengel. Vision-and-language navigation: Interpreting visually-grounded navigation instructions in real environments. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3674–3683, 2018. 8
- [4] Dhruv Batra, Angel X Chang, Sonia Chernova, Andrew J Davison, Jia Deng, Vladlen Koltun, Sergey Levine, Jitendra Malik, Igor Mordatch, Roozbeh Mottaghi, et al. Rearrangement: A challenge for embodied ai. *arXiv preprint arXiv:2011.01975*, 2020. 8
- [5] Samarth Brahmabhatt, Chengcheng Tang, Christopher D Twigg, Charles C Kemp, and James Hays. Contactpose: A dataset of grasps with object contact and hand pose. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XIII 16*, pages 361–378. Springer, 2020. 8
- [6] Holger Caesar, Varun Bankiti, Alex H. Lang, Sourabh Vora, Venice Erin Liong, Qiang Xu, Anush Krishnan, Yuxin Pan, Giancarlo Baldan, and Oscar Beijbom. nuscenes: A multimodal dataset for autonomous driving. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 11618–11628, 2019. 8
- [7] Angel Chang, Angela Dai, Thomas Funkhouser, Maciej Halber, Matthias Niessner, Manolis Savva, Shuran Song, Andy Zeng, and Yinda Zhang. Matterport3d: Learning from rgb-d data in indoor environments. *arXiv preprint arXiv:1709.06158*, 2017. 8
- [8] Angel X Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, et al. Shapenet: An information-rich 3d model repository. *arXiv preprint arXiv:1512.03012*, 2015. 3
- [9] Enric Corona, Albert Pumarola, Guillem Alenya, Francesc Moreno-Noguer, and Grégory Rogez. Ganhand: Predicting human grasp affordances in multi-object scenes. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5031–5041, 2020. 8
- [10] Abhishek Das, Samyak Datta, Georgia Gkioxari, Stefan Lee, Devi Parikh, and Dhruv Batra. Embodied question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1–10, 2018. 8
- [11] Abhishek Das, Federico Carnevale, Hamza Merzic, Laura Rimell, Rosalia Schneider, Josh Abramson, Alden Hung, Arun Ahuja, Stephen Clark, Gregory Wayne, et al. Probing emergent semantics in predictive agents via question answering. *arXiv preprint arXiv:2006.01016*, 2020. 8
- [12] Sudeep Dasari, Frederik Ebert, Stephen Tian, Suraj Nair, Bernadette Bucher, Karl Schmeckpeper, Siddharth Singh, Sergey Levine, and Chelsea Finn. Robonet: Large-scale multi-robot learning. *ArXiv*, abs/1910.11215, 2019. 8
- [13] Jiafei Duan, Wentao Yuan, Wilbert Pumacay, Yi Ru Wang, Kiana Ehsani, Dieter Fox, and Ranjay Krishna. Manipulate-anything: Automating real-world robots using vision-language models. *arXiv preprint arXiv:2406.18915*, 2024. 8
- [14] Kiana Ehsani, Tanmay Gupta, Rose Hendrix, Jordi Salvador, Luca Weihs, Kuo-Hao Zeng, Kunal Pratap Singh, Yejin Kim, Winson Han, Alvaro Herrasti, et al. Spoc: Imitating shortest paths in simulation enables effective navigation and manipulation in the real world. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16238–16250, 2024. 8
- [15] Zicong Fan, Omid Taheri, Dimitrios Tzionas, Muhammed Kocabas, Manuel Kaufmann, Michael J Black, and Otmar Hilliges. Arctic: A dataset for dexterous bimanual hand-object manipulation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12943–12954, 2023. 8
- [16] Zipeng Fu, Xuxin Cheng, and Deepak Pathak. Deep whole-body control: learning a unified policy for manipulation and locomotion. In *Conference on Robot Learning*, pages 138–149. PMLR, 2023. 8

- [17] Zipeng Fu, Tony Z Zhao, and Chelsea Finn. Mobile aloha: Learning bimanual mobile manipulation with low-cost whole-body teleoperation. *arXiv preprint arXiv:2401.02117*, 2024. 8
- [18] Xianqiang Gao, Pingrui Zhang, Delin Qu, Dong Wang, Zhigang Wang, Yan Ding, Bin Zhao, and Xuelong Li. Learning 2d invariant affordance knowledge for 3d affordance grounding. *arXiv preprint arXiv:2408.13024*, 2024. 8
- [19] Andreas Geiger, Philip Lenz, Christoph Stiller, and Raquel Urtasun. Vision meets robotics: The kitti dataset. *The International Journal of Robotics Research*, 32:1231 – 1237, 2013. 8
- [20] Kristen Grauman, Andrew Westbury, Eugene Byrne, Zachary Chavis, Antonino Furnari, Rohit Girdhar, Jackson Hamburger, Hao Jiang, Miao Liu, Xingyu Liu, et al. Ego4d: Around the world in 3,000 hours of egocentric video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18995–19012, 2022. 8
- [21] Huy Ha, Pete Florence, and Shuran Song. Scaling up and distilling down: Language-guided robot skill acquisition. In *Conference on Robot Learning*, pages 3766–3777. PMLR, 2023. 8
- [22] Wenlong Huang, Chen Wang, Ruohan Zhang, Yunzhu Li, Jiajun Wu, and Li Fei-Fei. Voxposer: Composable 3d value maps for robotic manipulation with language models. *ArXiv*, abs/2307.05973, 2023. 4
- [23] Juntao Jian, Xiuping Liu, Manyi Li, Ruizhen Hu, and Jian Liu. Affordpose: A large-scale dataset of hand-object interactions with affordance-driven hand pose. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14713–14724, 2023. 8
- [24] Yunfan Jiang, Ruohan Zhang, Josiah Wong, Chen Wang, Yanjie Ze, Hang Yin, Cem Gokmen, Shuran Song, Jiajun Wu, and Fei-Fei Li. Behavior robot suite: Streamlining real-world whole-body manipulation for everyday household activities. 2025. 8
- [25] Mukul Khanna, Yongsan Mao, Hanxiao Jiang, Sanjay Haresh, Brennan Shacklett, Dhruv Batra, Alexander Clegg, Eric Undersander, Angel X Chang, and Manolis Savva. Habitat synthetic scenes dataset (hssd-200): An analysis of 3d scene scale and realism tradeoffs for objectgoal navigation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16384–16393, 2024. 8
- [26] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization, 2017. 15
- [27] Jacob Krantz, Erik Wijmans, Arjun Majumdar, Dhruv Batra, and Stefan Lee. Beyond the nav-graph: Vision-and-language navigation in continuous environments. In *ECCV*, 2020. 8
- [28] Alexander Ku, Peter Anderson, Roma Patel, Eugene Ie, and Jason Baldridge. Room-across-room: Multilingual vision-and-language navigation with dense spatiotemporal grounding. *arXiv preprint arXiv:2010.07954*, 2020. 8
- [29] Sergey Levine, Peter Pastor, Alex Krizhevsky, and Deirdre Quillen. Learning hand-eye coordination for robotic grasping with deep learning and large-scale data collection. *The International Journal of Robotics Research*, 37:421 – 436, 2016. 8
- [30] Chengshu Li, Fei Xia, Roberto Martin-Martin, and Silvio Savarese. Hrl4in: Hierarchical reinforcement learning for interactive navigation with mobile manipulators. In *Conference on Robot Learning*, pages 603–616. PMLR, 2020. 8
- [31] Xunsong Li, Pengzhan Sun, Yangcen Liu, Lixin Duan, and Wen Li. Simultaneous detection and interaction reasoning for object-centric action recognition. *ArXiv*, abs/2404.11903, 2024. 8
- [32] Kehui Liu, Zixin Tang, Dong Wang, Zhigang Wang, Bin Zhao, and Xuelong Li. Coherent: Collaboration of heterogeneous multi-robot system with large language models. *arXiv preprint arXiv:2409.15146*, 2024. 8
- [33] Peiqi Liu, Yaswanth Orru, Jay Vakil, Chris Paxton, Nur Muhammad Mahi Shafiullah, and Lerrel Pinto. Ok-robot: What really matters in integrating open-knowledge models for robotics. *arXiv preprint arXiv:2401.12202*, 2024. 8
- [34] Yangcen Liu, Ziyi Liu, Yuanhao Zhai, Wen Li, David Doerman, and Junsong Yuan. Stat: Towards generalizable temporal action localization. *arXiv preprint arXiv:2404.13311*, 2024. 8
- [35] Ajay Mandlekar, Yuke Zhu, Animesh Garg, Jonathan Booher, Max Spero, Albert Tung, Julian Gao, John Emmons, Anchit Gupta, Emre Orbay, Silvio Savarese, and Li Fei-Fei. Roboturk: A crowdsourcing platform for robotic skill learning through imitation. In *Conference on Robot Learning*, 2018. 8
- [36] Karan Mirakhor, Sourav Ghosh, Dipanjan Das, and Brojeshwar Bhowmick. Task planning for visual room rearrangement under partial observability. In *The Twelfth International Conference on Learning Representations*, 2024. 8
- [37] Mayank Mittal, David Hoeller, Farbod Farshidian, Marco Hutter, and Animesh Garg. Articulated object interaction in unknown scenes with whole-body mobile manipulation. In *2022 IEEE/RSJ international conference on intelligent robots and systems (IROS)*, pages 1647–1654. IEEE, 2022. 8
- [38] Kaichun Mo, Shilin Zhu, Angel X. Chang, Li Yi, Subarna Tripathi, Leonidas J. Guibas, and Hao Su. PartNet: A large-scale benchmark for fine-grained and hierarchical part-level 3D object understanding. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 3
- [39] Lerrel Pinto and Abhinav Kumar Gupta. Supersizing self-supervision: Learning to grasp from 50k tries and 700 robot hours. *2016 IEEE International Conference on Robotics and Automation (ICRA)*, pages 3406–3413, 2015. 8
- [40] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 652–660, 2017. 15
- [41] Charles Ruizhongtai Qi, Li Yi, Hao Su, and Leonidas J Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. *Advances in neural information processing systems*, 30, 2017. 4, 6, 15
- [42] Charles R. Qi, Or Litany, Kaiming He, and Leonidas J. Guibas. Deep hough voting for 3d object detection in point clouds. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019. 6, 15

- [43] Yuankai Qi, Qi Wu, Peter Anderson, Xin Wang, William Yang Wang, Chunhua Shen, and Anton van den Hengel. Reverie: Remote embodied visual referring expression in real indoor environments. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9982–9991, 2020. 8
- [44] Dicong Qiu, Wenzong Ma, Zhenfu Pan, Hui Xiong, and Junwei Liang. Open-vocabulary mobile manipulation in unseen dynamic environments with 3d semantic maps. *arXiv preprint arXiv:2406.18115*, 2024. 2, 8
- [45] Santhosh K Ramakrishnan, Aaron Gokaslan, Erik Wijmans, Oleksandr Maksymets, Alex Clegg, John Turner, Eric Undersander, Wojciech Galuba, Andrew Westbury, Angel X Chang, et al. Habitat-matterport 3d dataset (hm3d): 1000 large-scale 3d environments for embodied ai. *arXiv preprint arXiv:2109.08238*, 2021. 8
- [46] Tianhe Ren, Shilong Liu, Ailing Zeng, Jing Lin, Kunchang Li, He Cao, Jiayu Chen, Xinyu Huang, Yukang Chen, Feng Yan, Zhaoyang Zeng, Hao Zhang, Feng Li, Jie Yang, Hongyang Li, Qing Jiang, and Lei Zhang. Grounded sam: Assembling open-world models for diverse visual tasks. *ArXiv*, abs/2401.14159, 2024. 7
- [47] Beichen Shao, Nieqing Cao, Yan Ding, Xingchen Wang, Fuqiang Gu, and Chao Chen. Moma-pos: An efficient object-kinematic-aware base placement optimization framework for mobile manipulation. *arXiv preprint arXiv:2403.19940*, 2024. 2
- [48] Pratyusha Sharma, Lekha Mohan, Lerrel Pinto, and Abhinav Kumar Gupta. Multiple interactions made easy (mime): Large scale demonstrations data for imitation. In *Conference on Robot Learning*, 2018. 8
- [49] Sanjana Srivastava, Chengshu Li, Michael Lingelbach, Roberto Martín-Martín, Fei Xia, Kent Elliott Vainio, Zheng Lian, Cem Gokmen, Shyamal Buch, Karen Liu, et al. Behavior: Benchmark for everyday household activities in virtual, interactive, and ecological environments. In *Conference on robot learning*, pages 477–490. PMLR, 2022. 8
- [50] Mike Stilman and James J Kuffner. Navigation among movable obstacles: Real-time reasoning in complex environments. *International Journal of Humanoid Robotics*, 2(04): 479–503, 2005. 8
- [51] Pei Sun, Henrik Kretschmar, Xerxes Dotiwalla, Aurelien Chouard, Vijaysai Patnaik, Paul Tsui, James Guo, Yin Zhou, Yuning Chai, Benjamin Caine, Vijay Vasudevan, Wei Han, Jiquan Ngiam, Hang Zhao, Aleksei Timofeev, Scott M. Etinger, Maxim Krivokon, Amy Gao, Aditya Joshi, Yu Zhang, Jonathon Shlens, Zhifeng Chen, and Dragomir Anguelov. Scalability in perception for autonomous driving: Waymo open dataset. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2443–2451, 2019. 8
- [52] Omid Taheri, Nima Ghorbani, Michael J Black, and Dimitrios Tzionas. Grab: A dataset of whole-body human grasping of objects. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part IV 16*, pages 581–600. Springer, 2020. 8
- [53] Yiwen Tang, Ivan Tang, Eric Zhang, and Ray Gu. Point-peft: Parameter-efficient fine-tuning for 3d pre-trained models. *ArXiv*, abs/2310.03059, 2023. 8
- [54] Yiwen Tang, Jiaming Liu, Dong Wang, Zhigang Wang, Shanghang Zhang, Bin Zhao, and Xuelong Li. Any2point: Empowering any-modality large models for efficient 3d understanding. *ArXiv*, abs/2404.07989, 2024. 8
- [55] Jur Van Den Berg, Sachin Patil, Jason Sewall, Dinesh Manocha, and Ming Lin. Interactive navigation of multiple agents in crowded environments. In *Proceedings of the 2008 symposium on Interactive 3D graphics and games*, pages 139–147, 2008. 8
- [56] A Vaswani. Attention is all you need. *Advances in Neural Information Processing Systems*, 2017. 5
- [57] Weikang Wan, Haoran Geng, Yun Liu, Zikang Shan, Yaodong Yang, Li Yi, and He Wang. Unidexgrasp++: Improving dexterous grasping policy learning via geometry-aware curriculum and iterative generalist-specialist learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3891–3902, 2023. 8
- [58] Huaxiaoyue Wang, Kushal Kedia, Juntao Ren, Rahma Abdullah, Atiksh Bhardwaj, Angela Chao, Kelly Y Chen, Nathaniel Chin, Prithwish Dan, Xinyi Fan, et al. Mosaic: A modular system for assistive and interactive cooking. *arXiv preprint arXiv:2402.18796*, 2024. 8
- [59] Lirui Wang, Yiyang Ling, Zhecheng Yuan, Mohit Shridhar, Chen Bao, Yuzhe Qin, Bailin Wang, Huazhe Xu, and Xiaolong Wang. Gensim: Generating robotic simulation tasks via large language models. In *Arxiv*, 2023. 8
- [60] Lirui Wang, Xinlei Chen, Jialiang Zhao, and Kaiming He. Scaling proprioceptive-visual learning with heterogeneous pre-trained transformers. *arXiv preprint arXiv:2409.20537*, 2024. 2
- [61] Maozhen Wang, Rui Luo, Aykut Özgün Önel, and Taşkın Padir. Affordance-based mobile robot navigation among movable obstacles. In *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 2734–2740. IEEE, 2020. 8
- [62] Ruicheng Wang, Jialiang Zhang, Jiayi Chen, Yinzheng Xu, Puhao Li, Tengyu Liu, and He Wang. Dexgraspnet: A large-scale robotic dexterous grasp dataset for general objects based on simulation. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 11359–11366. IEEE, 2023. 8
- [63] Xiaohan Wang, Yuehu Liu, Xinhang Song, Beibei Wang, and Shuqiang Jiang. Camp: Causal multi-policy planning for interactive navigation in multi-room scenes. *Advances in Neural Information Processing Systems*, 36, 2024. 8
- [64] Yufei Wang, Zhou Xian, Feng Chen, Tsun-Hsuan Wang, Yian Wang, Katerina Fragkiadaki, Zackory Erickson, David Held, and Chuang Gan. Robogen: Towards unleashing infinite data for automated robot learning via generative simulation. *arXiv preprint arXiv:2311.01455*, 2023. 8
- [65] Luca Weihs, Matt Deitke, Aniruddha Kembhavi, and Roozbeh Mottaghi. Visual room rearrangement. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5922–5931, 2021. 8
- [66] Erik Wijmans, Samyak Datta, Oleksandr Maksymets, Abhishek Das, Georgia Gkioxari, Stefan Lee, Irfan Essa, Devi

- Parikh, and Dhruv Batra. Embodied question answering in photorealistic environments with point cloud perception. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6659–6668, 2019. 8
- [67] Josiah Wong, Albert Tung, Andrey Kurenkov, Ajay Mandlekar, Li Fei-Fei, Silvio Savarese, and Roberto Martín-Martín. Error-aware imitation learning from teleoperation data for mobile manipulation. In *Conference on Robot Learning*, pages 1367–1378. PMLR, 2022. 8
- [68] Bohan Wu, Roberto Martín-Martín, and Li Fei-Fei. Member: Tackling long-horizon mobile manipulation via factorized domain transfer. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 11690–11697. IEEE, 2023. 8
- [69] Fei Xia, Amir R Zamir, Zhiyang He, Alexander Sax, Jitendra Malik, and Silvio Savarese. Gibson env: Real-world perception for embodied agents. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 9068–9079, 2018. 8
- [70] Fei Xia, William B Shen, Chengshu Li, Priya Kasimbeg, Micael Edmond Tchampi, Alexander Toshev, Roberto Martín-Martín, and Silvio Savarese. Interactive gibbon benchmark: A benchmark for interactive navigation in cluttered environments. *IEEE Robotics and Automation Letters*, 5(2):713–720, 2020. 8
- [71] Fanbo Xiang, Yuzhe Qin, Kaichun Mo, Yikuan Xia, Hao Zhu, Fangchen Liu, Minghua Liu, Hanxiao Jiang, Yifu Yuan, He Wang, Li Yi, Angel X. Chang, Leonidas J. Guibas, and Hao Su. SAPIEN: A simulated part-based interactive environment. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 3
- [72] Jiannan Xiang, Tianhua Tao, Yi Gu, Tianmin Shu, Zirui Wang, Zichao Yang, and Zhiting Hu. Language models meet world models: Embodied experiences enhance language models. In *Advances in Neural Information Processing Systems*, pages 75392–75412. Curran Associates, Inc., 2023. 2
- [73] Yinzhen Xu, Weikang Wan, Jialiang Zhang, Haoran Liu, Zikang Shan, Hao Shen, Ruicheng Wang, Haoran Geng, Yijia Weng, Jiayi Chen, et al. Unidexgrasp: Universal robotic dexterous grasping via learning diverse proposal generation and goal-conditioned policy. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4737–4746, 2023. 8
- [74] Kui Yang, Nieqing Cao, Yan Ding, and Chao Chen. Bestman: A modular mobile manipulator platform for embodied ai with unified simulation-hardware apis. 2024. 2, 4, 13
- [75] Lixin Yang, Kailin Li, Xinyu Zhan, Fei Wu, Anran Xu, Liu Liu, and Cewu Lu. Oakink: A large-scale knowledge repository for understanding hand-object interaction. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 20953–20962, 2022. 8
- [76] Lihe Yang, Bingyi Kang, Zilong Huang, Zhen Zhao, Xiaogang Xu, Jiashi Feng, and Hengshuang Zhao. Depth anything v2. *ArXiv*, abs/2406.09414, 2024. 8
- [77] Sriram Yenamandra, Arun Ramachandran, Karmesh Yadav, Austin Wang, Mukul Khanna, Theophile Gervet, Tsung-Yen Yang, Vidhi Jain, Alexander William Clegg, John Turner, et al. Homerobot: Open-vocabulary mobile manipulation. *arXiv preprint arXiv:2306.11565*, 2023. 8
- [78] Naoki Yokoyama, Alex Clegg, Joanne Truong, Eric Under-sander, Tsung-Yen Yang, Sergio Arnaud, Sehoon Ha, Dhruv Batra, and Akshara Rai. Asc: Adaptive skill coordination for robotic mobile manipulation. *IEEE Robotics and Automation Letters*, 9(1):779–786, 2023. 8
- [79] Licheng Yu, Xinlei Chen, Georgia Gkioxari, Mohit Bansal, Tamara L Berg, and Dhruv Batra. Multi-target embodied question answering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6309–6318, 2019. 8
- [80] Kuo-Hao Zeng, Luca Weihs, Ali Farhadi, and Roozbeh Motaghi. Pushing it out of the way: Interactive visual navigation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9868–9877, 2021. 8
- [81] Xinyu Zhan, Lixin Yang, Yifei Zhao, Kangrui Mao, Hanlin Xu, Zenan Lin, Kailin Li, and Cewu Lu. Oakink2: A dataset of bimanual hands-object manipulation in complex task completion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 445–456, 2024. 8
- [82] Yu Zhang, Kehai Chen, Xuefeng Bai, Zhao Kang, Quanjiang Guo, and Min Zhang. Question-guided knowledge graph rescoring and injection for knowledge graph question answering. In *Conference on Empirical Methods in Natural Language Processing*, 2024. 8
- [83] Zaiwei Zhang, Bo Sun, Haitao Yang, and Qi-Xing Huang. H3dnet: 3d object detection using hybrid geometric primitives. In *European Conference on Computer Vision*, 2020. 6, 15
- [84] Zhaxizhuoma, Kehui Liu, Chuyue Guan, Zhongjie Jia, Ziniu Wu, Xin Liu, Tianyu Wang, Shuai Liang, Pengan Chen, Pingrui Zhang, Haoming Song, Delin Qu, Dong Wang, Zhigang Wang, Nieqing Cao, Yan Ding, Bin Zhao, and Xuelong Li. Fastumi: A scalable and hardware-independent universal manipulation interface with dataset, 2025. 8
- [85] Peiyuan Zhi, Zhiyuan Zhang, Muzhi Han, Zeyu Zhang, Zhitian Li, Ziyuan Jiao, Baoxiong Jia, and Siyuan Huang. Closed-loop open-vocabulary mobile manipulation with gpt-4v. *arXiv preprint arXiv:2404.10220*, 2024. 8
- [86] Peiyuan Zhi, Zhiyuan Zhang, Muzhi Han, Zeyu Zhang, Zhitian Li, Ziyuan Jiao, Baoxiong Jia, and Siyuan Huang. Closed-loop open-vocabulary mobile manipulation with gpt-4v, 2024. 2
- [87] Gengze Zhou, Yicong Hong, Zun Wang, Xin Eric Wang, and Qi Wu. Navgpt-2: Unleashing navigational reasoning capability for large vision-language models. *arXiv preprint arXiv:2407.12366*, 2024. 2
- [88] Jun Zhu, Zihao Du, Haotian Xu, Fengbo Lan, Zilong Zheng, Bo Ma, Shengjie Wang, and Tao Zhang. Navi2gaze: Leveraging foundation models for navigation and target gazing. *arXiv preprint arXiv:2407.09053*, 2024. 2

# MoMa-Kitchen: A 100K+ Benchmark for Affordance-Grounded Last-Mile Navigation in Mobile Manipulation

## Supplementary Material

This supplementary material extends our main study by providing additional details and data to improve the reproducibility of our MoMa-Kitchen method. It includes further evaluations and a range of qualitative results for NavAff, which reinforce the conclusions drawn in the primary paper. Additionally, we offer some affordance collection videos in the accompanying zip file.

▷ **Sec. 1:** Describes the hierarchical structure of the dataset, including scenes, configurations, and episodes, with detailed information on the generation process, target objects, and simulation settings.

▷ **Sec. 2:** Provides an in-depth explanation of the evaluation metrics used, training configurations, and the baseline models compared in our study.

▷ **Sec. 3:** Presents additional visualizations of predictions, further performance evaluations, ablation results regarding the weight of the MSE loss, some data collection videos, and real-world demo video.

▷ **Sec. 4:** Discusses the limitations of our work and explores prospects for future research.

## 1. Dataset

### 1.1. Dataset Composition and Splits

Our MoMa-Kitchen is hierarchically organized into three levels: scenes, configurations, and episodes. Here, we provide a detailed description of each level.

A scene consists of randomly generated base furniture and layout, where certain articulated objects in the base furniture (e.g., microwaves, oven counters) serve as potential target objects for robotic arm manipulation. To ensure scene diversity, we randomly sample furniture categories, arrangement sequences, and specific instances within categories during scene generation. The statistics of target object assets employed in MoMa-Kitchen are summarized in Tab. 5.

Within each scene, we randomly place a varying number (1-3) of rigid objects, which, together with the articulated objects, constitute the set of target objects. To increase scene complexity, we position obstacles around these target objects. Each unique combination of target objects and obstacles forms a configuration of the scene.

To facilitate first-person view data collection, we sample 10 views for each configuration using a camera mounted on the robotic arm. Each view generates one episode, and the view selection follows two principles: (1) views are ran-

Rigid-Cats	All	Bottle	Pot	Fruit	Medicine Bottle	Vegetable
Rigid-Num	65	6	7	11	8	33
Articulated-Cats	All	Faucet	Microwave	Cabinet	Dishwasher	Oven Counter
Articulated-Num	48	11	7	20	1	9
Obstacle-Cats	All	Chair	Trolley	Bin	Table	Cart
Obstacle-Num	24	1	8	1	10	4

Table 5. **Statistics of target object assets employed in MoMa-Kitchen.** Distribution of rigid, articulated, and obstacle objects across different categories, showing the number of instances per category used in our dataset configurations.

domly initialized around the target, and (2) views must encompass both the target object and the surrounding floor area.

In total, our MoMa-Kitchen comprises 569 scenes, 14, 155 configurations, and 127, 343 episodes, representing a comprehensive collection of mobile manipulation scenarios.

### 1.2. Details on Simulation

We build our MoMa-Kitchen based on the BestMan [74] simulation environment, maintaining consistent simulation parameters across all scenes and interaction trials. The detailed configuration of our simulation setup is specified below:

- **RGBD rendering.** We render RGB images and depth maps using the BestMan interface. For comprehensive first-person view sampling, we position the camera at varying locations relative to the target object. Specifically, the camera is placed either to the left or right with a lateral offset ranging from 0.0 to 1.5 meters, while the forward distance is sampled between 1.5 and 3.8 meters. The camera orientation is consistently directed toward the target object to ensure optimal coverage of both the target and the surrounding floor area. These sampling ranges were empirically determined to maximize viewpoint diversity while maintaining scene relevance.
- **3D point cloud.** We back-project the depth image into a 3D point cloud using the camera’s intrinsic parameters. Subsequently, we filter out points with z-values below 0.02 meters to obtain the floor point cloud.
- **Target objects sampling.** In each scene, we randomly position 1-3 rigid objects in addition to the pre-existing articulated objects from scene generation. These objects

collectively form our set of target objects, all of which are placed on kitchen countertops. To increase scene complexity, we randomly place 1-3 obstacles within the semi-circular region in front of each target object.

- **Interaction Trail.** To collect discrete navigation affordance values, we systematically sample robot positions within a semicircular region around the target object, with the radius set to the maximum reach of the robot arm. At each position, spaced at 10 cm intervals along both x and y axes, the robot attempts to either grasp (for parallel grippers) or suction (for vacuum grippers) the target object.

### 1.3. Additional Visualization Results

We present additional visualization examples from MoMa-Kitchen, including object assets, scene configurations, and affordance maps.

#### 1.3.1. Object Assets

We showcase a diverse set of object assets in Fig. 8. These assets serve as manipulation targets, environmental elements, or obstacles in our generated scenes.

#### 1.3.2. Scene Configurations

We illustrate a comprehensive set of scene configurations in Fig. 9. These examples highlight the complexity and diversity of our generated environments, reflecting real-world manipulation scenarios.

#### 1.3.3. Affordance Map Examples

We obtain sparse discrete navigation affordance values by systematically sampling robot positions and evaluating object interactions in the simulator. These sparse affordance values are then interpolated using a Gaussian-weighted k-nearest neighbor algorithm to generate continuous and dense navigation affordance maps. The comparison between sparse samples and interpolated dense maps is illustrated in Fig. 7.

## 2. Additional Implementation Details

### 2.1. Evaluation Metrics

In this section, we provide a detailed explanation of the evaluation metrics used in our study. To comprehensively evaluate our method, we adopt five metrics: Root Mean Squared Error (**RMSE**), Logarithmic Mean Squared Error (**logMSE**), Pearson Correlation Coefficient (**PCC**), Cosine Similarity (**SIM**), and Continuous Intersection over Union (**cIoU**). Below, we describe each metric, its calculation formula, and its relevance to our task.

- **RMSE:** RMSE measures the numerical alignment between predicted and ground truth values. It penalizes larger errors through squared differences, as defined be-

low:

$$\text{RMSE} = \sqrt{\frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2}, \quad (11)$$

where  $y_i$  and  $\hat{y}_i$  are ground truth and predicted values, respectively, and  $N$  is the total number of elements. RMSE reintroduces the original scale of predictions, offering interpretability while highlighting large errors. In our study, RMSE assesses navigation affordance predictions by ensuring precise positioning and minimizing major errors in complex environments.

- **logMSE:** logMSE evaluates the relative differences between predictions and ground truth, reducing the impact of large outliers. It is calculated as:

$$\text{logMSE} = \frac{1}{N} \sum_{i=1}^N (\log(1 + y_i) - \log(1 + \hat{y}_i))^2, \quad (12)$$

By focusing on proportional consistency, logMSE smooths out outliers and highlights relative accuracy. In this study, it evaluates the model’s ability to capture balanced affordance patterns across both low and high value regions.

- **PCC:** PCC quantifies the linear relationship between predicted and ground truth patterns, independent of their magnitudes:

$$\text{PCC} = \frac{\sum_{i=1}^N (y_i - \bar{y})(\hat{y}_i - \bar{\hat{y}})}{\sqrt{\sum_{i=1}^N (y_i - \bar{y})^2 \sum_{i=1}^N (\hat{y}_i - \bar{\hat{y}})^2}}, \quad (13)$$

where  $\bar{y}$  and  $\bar{\hat{y}}$  are the means of  $y$  and  $\hat{y}$ . PCC highlights pattern consistency, making it ideal for evaluating spatial distributions in affordance maps. A high PCC reflects accurate predictions of affordance trends, which is essential for precise navigation.

- **SIM:** SIM evaluates the alignment of relative spatial patterns between predictions and ground truth. It is calculated as:

$$\text{SIM} = \frac{1}{N} \sum_{i=1}^N \frac{y_i \cdot \hat{y}_i}{\|y_i\| \|\hat{y}_i\|} \quad (14)$$

SIM with higher values indicating better structural alignment. Unlike PCC, which evaluates overall pattern trends, SIM focuses on the overlap of affordance regions, making it particularly effective for assessing spatial alignment. In our study, SIM ensures that predicted affordance maps accurately capture the structural properties of ground truth.

### 2.2. Training Details

To ensure a fair comparison, we train our model and all baseline methods under consistent settings. The implementations for both our methods and the baselines are developed

using PyTorch. All models are trained on a single NVIDIA A100 GPU with a batch size of 64 for 6 epochs, completing the entire process in approximately 8 hours. We utilize the Adam optimizer [26] with betas configured as 0.9 and 0.999. The learning rate is initialized at  $8e-4$  and follows a cosine decay schedule.

### 2.3. Compared Baselines

Since our work is the first to propose a benchmark for navigation affordance grounding, there are no existing methods that directly address this task. Therefore, we adapt several classical methods commonly used for feature extraction and 3D object detection on point clouds for comparison. Specifically, we include PointNet++[41], a foundational model for point cloud feature extraction, VoteNet[42], a pioneering method for 3D object proposal generation, and H3DNet [83], which enhances object detection with hierarchical features. To ensure a fair and comprehensive comparison, we adapt the official implementations of these methods to our MoMa-Kitchen dataset. Specifically, we reimplement their architectures and fine-tune them for the navigation affordance grounding task, conducting training and evaluation under identical experimental settings. This allows us to systematically assess their performance against our proposed NavAff.

**PointNet++.** PointNet++ extends the original PointNet [40] framework by introducing hierarchical feature learning for point cloud processing. This method divides the input point cloud into overlapping regions using a sampling and grouping strategy, applying PointNet locally to extract features, and aggregating them hierarchically. PointNet++ is widely used for tasks such as segmentation and classification in 3D point clouds. In our adaptation, the point cloud data is passed through an encoder to extract features, which are then decoded to predict navigation affordance.

**VoteNet.** VoteNet introduces a deep Hough voting framework for 3D object detection in point clouds. The method employs a point-based network to generate votes for object centers, followed by an aggregation module that clusters votes to produce 3D object proposals. To adapt VoteNet for our benchmark, we removed the Vote Aggregation and Detection components originally used for bounding box regression, retaining the remaining modules to perform navigation affordance grounding. This adaptation ensures the network focuses on predicting affordance maps instead of object detection.

**H3DNet.** H3DNet proposes a Hierarchical 3D Detection Network that improves 3D object detection by leveraging multi-level geometric features. The network integrates instance-level and part-level features using a coarse-to-fine detection pipeline and introduces novel feature aggregation modules to enhance geometric reasoning. Similar to VoteNet, in our adaptation, we removed the bounding box

regression components while retaining the remaining modules to focus on navigation affordance grounding, enabling the network to predict affordance maps instead of object detection outcomes.

## 3. Additional Experimental Results

**Visualization of Predictions.** As shown in Fig. 10, we present the predicted navigation affordance results visualized within the dense global point cloud. Additionally, we provide a comparison with the ground truth affordance to highlight the model’s performance and alignment with the reference data.

**Detailed Evaluation.** Tab. 6 presents a comprehensive evaluation of various metrics within a single scene, offering an in-depth analysis of the model’s behavior and performance. Each scene comprises multiple episodes, each representing distinct configurations and challenges. By evaluating metrics across these episodes, we gain a finer understanding of the model’s ability to generalize under varying conditions.

**Impact of Weight Choices on Weight MSE Loss.** Fig. 11 illustrates that when the weight value is too small, the model experiences underfitting because the influence of the loss function weight is insufficient, preventing the model from effectively learning high-quality navigation affordance grounding. Conversely, when the weight value is too large, the class imbalance issue described earlier persists, which also limits the model’s performance. From the figure, it can be observed that when the weight value is set to 0.7, the Pearson Correlation Coefficient (PCC) reaches its peak. PCC measures the linear correlation between the predicted and ground truth values, effectively reflecting the model’s ability to capture the distribution patterns of navigation affordance. A high PCC value indicates a stronger correlation between the predicted trends and the ground truth distribution, which is particularly critical for navigation tasks in complex environments.

**Real world Experiment.** Please see the real-world captured video demo in the zip file.

**Affordance labeling visualization.** Please see the navigation affordance collection video in the zip file.

## 4. Discussion on Limitations and Future Work

While our work makes significant progress in addressing the “last mile” navigation challenge, we acknowledge several limitations and identify promising directions for future research:

- **Scene Diversity:** Although MoMa-Kitchen contains a large number of episodes, they are currently limited to kitchen environments. Future work should expand to other household scenarios such as living rooms, bedrooms, and bathrooms, which present different challenges

and spatial configurations.

- **Single-Task Focus:** The current approach focuses solely on reaching and grasping tasks. Future work should consider more complex manipulation sequences that require multiple positioning adjustments or different types of interactions (*e.g.*, pushing, pulling, or sliding objects).

**Future Directions**

- **Multi-Task Learning:** Future research could explore how navigation affordances vary across different manipulation tasks and develop models that can adapt their positioning strategies based on the intended manipulation action.
- **Online Adaptation:** Developing methods for online adjustment of affordance predictions based on real-time feedback during task execution could enhance robustness in dynamic environments.
- **Integration with LLMs:** While current LLM-based approaches have limitations, future work could explore hybrid approaches that combine our learned affordance models with LLM reasoning for more sophisticated task planning and execution.
- **Uncertainty Estimation:** Incorporating uncertainty estimation in affordance predictions could help robots make more informed decisions about positioning and potentially trigger replanning when necessary.

These limitations and future directions present exciting opportunities for extending our work and further advancing the field of mobile manipulation.

Table 6. **Evaluation results across individual scenes on MoMa-Kitchen.** Performance metrics evaluated separately for each scene in the dataset.

Scene ID	RMSE ↓	logMSE ↓	PCC ↑	SIM ↑
989172	0.283	0.0439	0.685	0.702
807952	0.269	0.0402	0.652	0.667
502334	0.284	0.0443	0.716	0.732
306938	0.282	0.0435	0.737	0.752
443958	0.232	0.0303	0.615	0.622
66171	0.297	0.0493	0.595	0.619
152285	0.242	0.0332	0.552	0.578
306168	0.223	0.0282	0.595	0.609
636942	0.283	0.0445	0.694	0.713
739657	0.264	0.0392	0.652	0.672
143853	0.291	0.0463	0.742	0.756
739202	0.236	0.0312	0.561	0.584
583009	0.302	0.0487	0.702	0.719
451797	0.301	0.0491	0.728	0.745
116280	0.166	0.0164	0.305	0.364
274269	0.274	0.0405	0.829	0.827
485779	0.298	0.0480	0.730	0.749
772552	0.299	0.0494	0.683	0.706
359363	0.213	0.0255	0.566	0.576
264325	0.292	0.0454	0.593	0.622
194561	0.276	0.0417	0.678	0.695
792629	0.292	0.0462	0.701	0.715
69567	0.273	0.0416	0.722	0.735
783647	0.288	0.0457	0.644	0.666
721930	0.219	0.0268	0.500	0.528
668061	0.266	0.0385	0.587	0.616
615543	0.275	0.0408	0.634	0.648
994972	0.285	0.0436	0.759	0.754
996121	0.266	0.0389	0.587	0.604
67534	0.278	0.0421	0.633	0.654
142672	0.270	0.0402	0.700	0.714
501160	0.276	0.0426	0.692	0.714
487375	0.227	0.0286	0.519	0.528
437964	0.294	0.0465	0.712	0.729
355986	0.299	0.0485	0.744	0.764
453020	0.285	0.0449	0.663	0.686
309033	0.296	0.0465	0.758	0.762
567413	0.225	0.0284	0.608	0.619
23304	0.248	0.0340	0.560	0.586
960190	0.247	0.0336	0.683	0.699
243997	0.288	0.0454	0.750	0.762
466622	0.265	0.0389	0.571	0.592
569661	0.278	0.0424	0.685	0.705
403556	0.195	0.0219	0.481	0.504
297024	0.289	0.0454	0.724	0.742
419493	0.219	0.0278	0.300	0.369
179882	0.213	0.0255	0.566	0.625
328786	0.256	0.0370	0.629	0.643
475545	0.186	0.0202	0.573	0.568
773991	0.278	0.0425	0.665	0.684



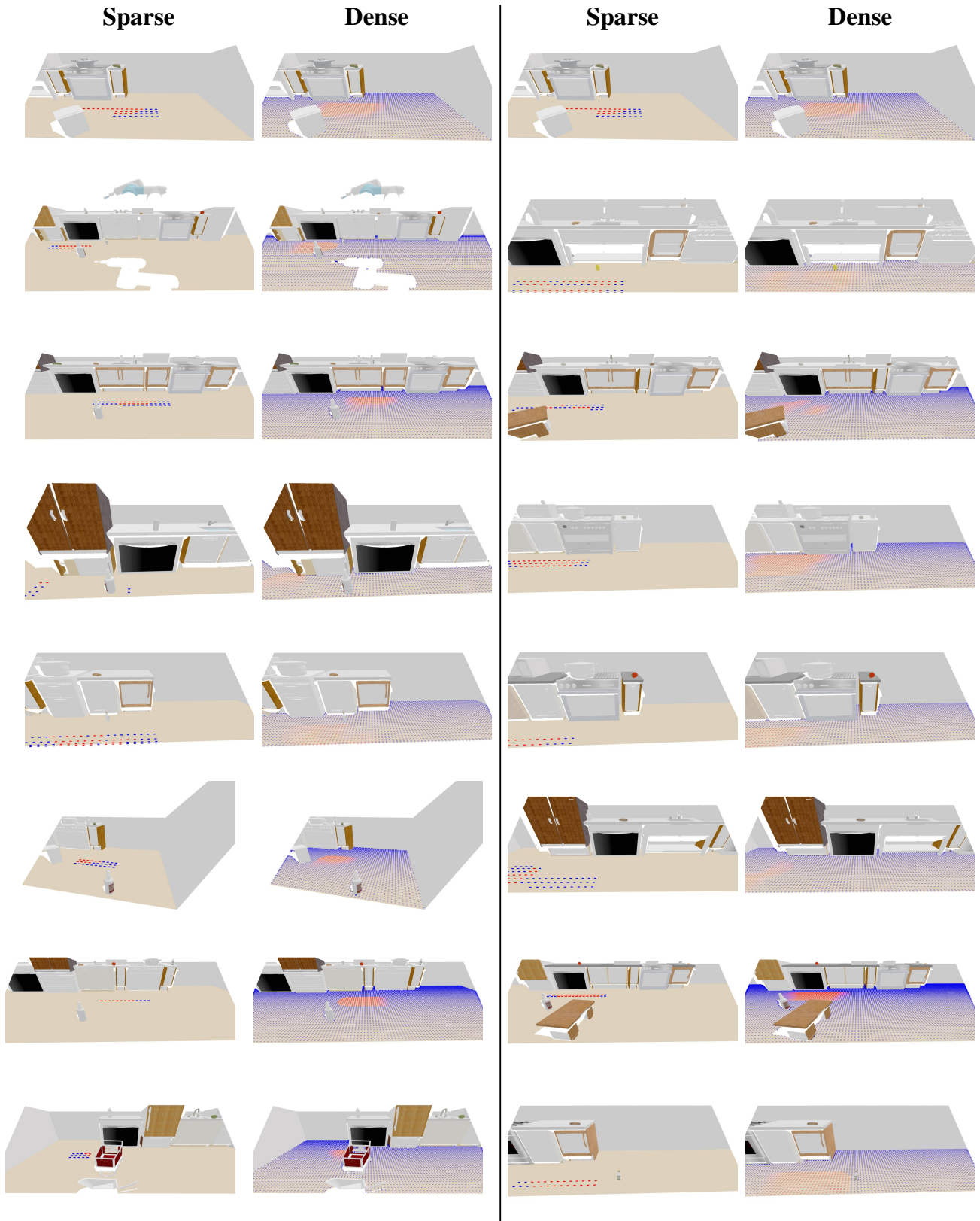


Figure 7. **Visualization of navigation affordance maps in MoMa-Kitchen.** Comparison between sparse affordance values collected through discrete robot-object interactions (left) and their corresponding dense maps generated via Gaussian-weighted k-nearest interpolation (right).



Figure 8. **Visualization of object assets in MoMa-Kitchen.** The collection includes diverse categories of objects commonly found in household environments, ranging from kitchenware and appliances to furniture and daily necessities.

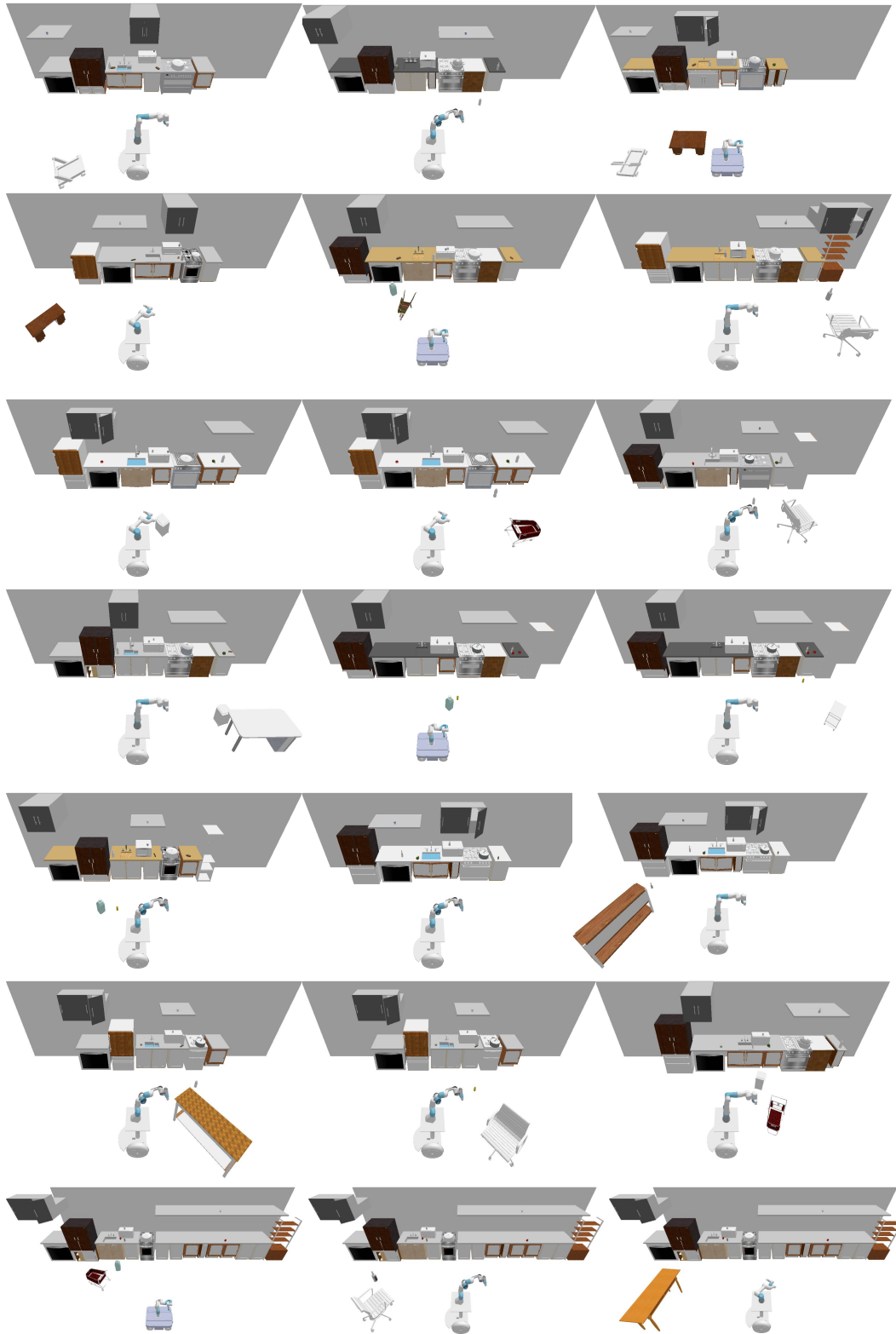


Figure 9. **Visualization of diverse scene configurations in MoMa-Kitchen.** Each example showcases different arrangements of base furniture, layouts, target objects, and obstacles, demonstrating the variety of manipulation scenarios.

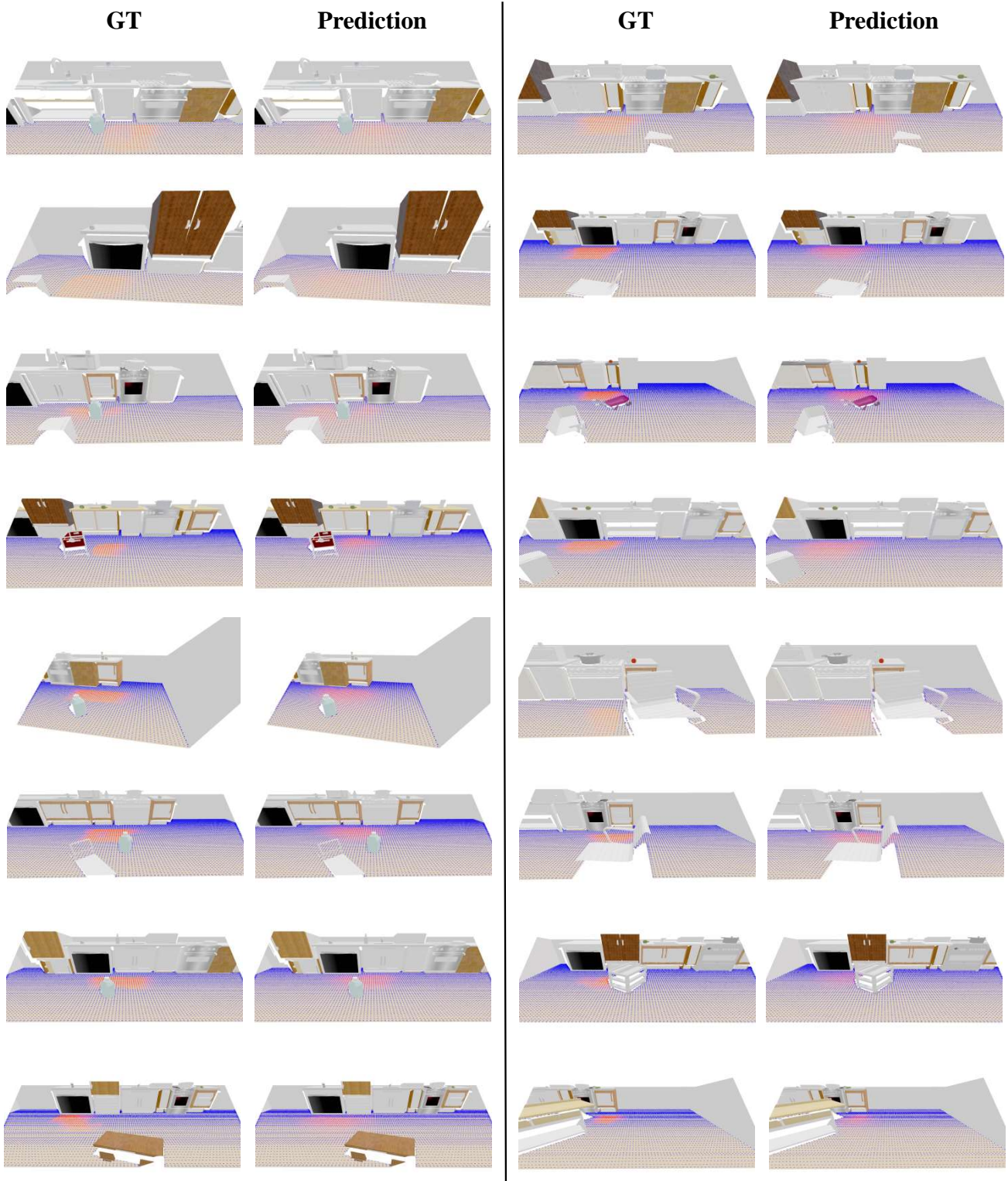


Figure 10. **Predicted vs. Ground Truth Navigation Affordance.** Comparison of the model’s predicted navigation affordance (right columns) and the ground truth affordance (left columns) visualized within dense global point clouds. The visualizations illustrate the spatial alignment and consistency of the predictions with the reference data across different scenes.



Figure 11. **Effect of Weight on MSE Loss and Evaluation Metrics.** Evaluation of the impact of different weight values in the Weighted MSE loss function on various metrics, including RMSE, LogMSE, PCC, and SIM.