# AnywhereVLA: Language-Conditioned Exploration and Mobile Manipulation

Konstantin Gubernatorov*, Artem Voronov*, Roman Voronov*, Sergei Pasynkov*,
Stepan Perminov, Ziang Guo, and Dzmitry Tsetserukou

*Abstract*— We address natural language pick-and-place in unseen, unpredictable indoor environments with AnywhereVLA, a modular framework for mobile manipulation. A user text prompt serves as an entry point and is parsed into a structured task graph that conditions classical SLAM with LiDAR and cameras, metric semantic mapping, and a task-aware frontier exploration policy. An approach planner then selects visibility and reachability aware pre grasp base poses. For interaction, a compact SmolVLA manipulation head is fine tuned on platform pick and place trajectories for the SO-101 by TheRobotStudio, grounding local visual context and sub-goals into grasp and place proposals. The full system runs fully onboard on consumer-level hardware, with Jetson Orin NX for perception and VLA and an Intel NUC for SLAM, exploration, and control, sustaining real-time operation. We evaluated AnywhereVLA in a multi-room lab under static scenes and normal human motion. In this setting, the system achieves a 46% overall task success rate while maintaining throughput on embedded compute. By combining a classical stack with a fine-tuned VLA manipulation, the system inherits the reliability of geometry-based navigation with the agility and task generalization of language-conditioned manipulation. All code, models, and datasets are open source and are available on the **project GitHub repository**.

## I. INTRODUCTION

Mobile manipulation is accelerating beyond limited indoor workcells towards large unstructured environments, in which robots need to explore unfamiliar cluttered spaces and physically interact with diverse objects and people. The execution of complex mobile manipulation tasks conditional on natural language instructions has gained attention in recent years in the field of service robotics [1]. Research on intelligent robotic systems in fields such as household service [2], [3], retail automation [4], [5], warehouse logistics [6], and manufacturing [7] has gained popularity, highlighting the importance of developing mobile manipulation solutions capable of operating in large-scale and open-plan indoor environments. Recent studies have increasingly focused on natural language processing to enable robots to interpret human instructions and facilitate intuitive task specification [8], positioning language-guided manipulation as a key approach for effective human-robot collaboration [9]. However, unifying language-based control, environment exploration, and manipulation in expansive environments presents a significant challenge [10].

*Denotes equal contribution.

All authors are with the Intelligent Space Robotics Laboratory, Center for Digital Engineering, Skolkovo Institute of Science and Technology, Moscow, Russia. {Konstantin.Gubernatorov, Artem.Voronov, Roman.Voronov, Sergei.Pasynkov, Stepan.Perminov, Ziang.Guo, D.Tsetserukou} @skoltech.ru

Vision-language-action (VLA) models show strong generalization in various mobile manipulation tasks [11], [12], enabling robots to perform complex operations integrating perception, language, and control. Despite these advancements, several critical limitations persist for end-to-end control use for mobile robots. Most VLA models are confined to specific tasks and have limited spatial awareness [11]–[13], restricting their operational scope to localized settings and impeding their ability to navigate or manipulate objects in unseen or occluded regions of larger indoor spaces.

Vision-Language Navigation (VLN) approaches [14] present an end-to-end VLM-based framework that navigates building-wide environments while manipulating previously unseen household objects. However, [14], as all VLN models, require instructions about the location of the target object within the environment, which is often impractical in dynamic or unexplored settings. In contrast, classical navigation stacks [15] provide robust solutions for mapping and environment exploration, enabling robots to systematically traverse and model unknown spaces. However, these traditional systems lack the advanced language comprehension and semantic reasoning capabilities necessary to interpret complex instructions or contextual cues [8], limiting their ability to perform goal-directed tasks that require understanding natural language or high-level semantic objectives.

In this paper we introduce AnywhereVLA, a novel modular architecture for large-scale indoor mobile manipulation, addressing the limitations of existing Vision-Language-Action (VLA) models that adapt pretrained vision-language models (VLMs) to enable natural language-driven perception and control, but are often constrained to room-scale environments due to high computational demands. Drawing from advancements in VLA paradigms, AnywhereVLA integrates the rich visual and linguistic knowledge encoded in VLMs—leveraged through co-fine-tuning on robotic trajectory data and Internet-scale vision-language tasks—with the robust traversability afforded by classical navigation stacks, representing robot actions as tokenized sequences to facilitate end-to-end control and emergent semantic reasoning.

AnywhereVLA is a pipeline for large-scale indoor mobile manipulation in unseen environments. As shown in Fig. 1, AnywhereVLA combines robust traversability of classical navigation algorithms and simultaneous localization and mapping (SLAM) with the generalizable scene understanding and task grounding of VLA models. Our pipeline translates high-level language instructions into low-level control commands by generating actions via VLA model for task-
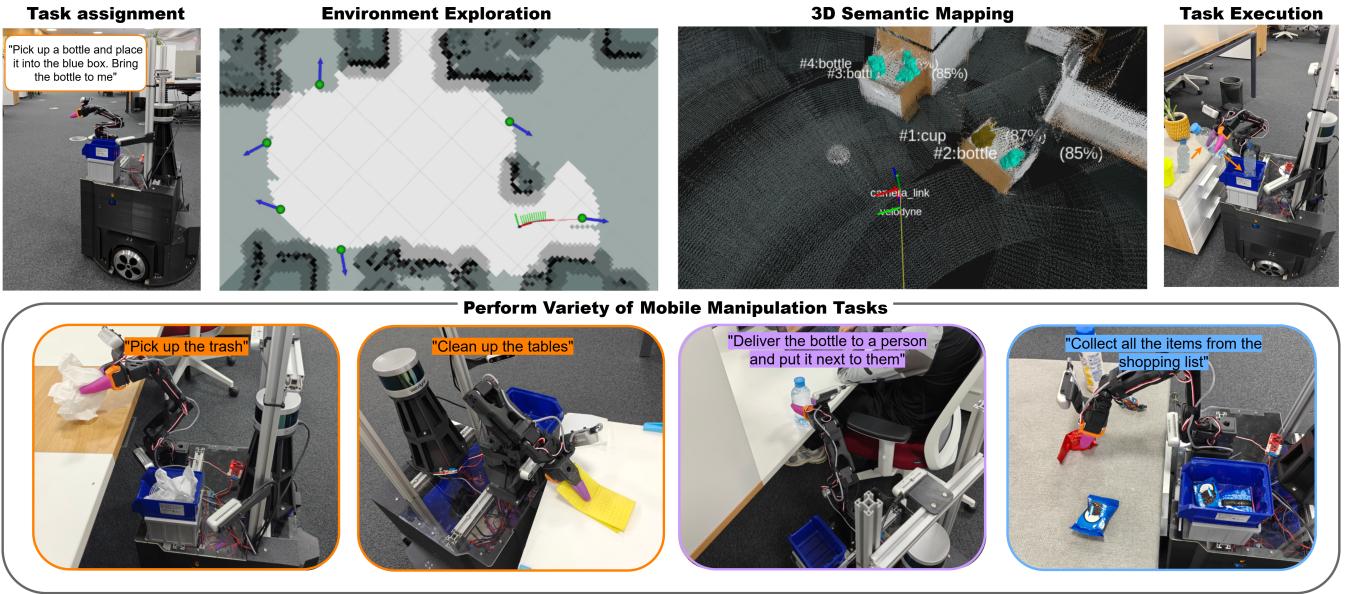
Fig. 1: AnywhereVLA is the modular architecture comprising VLA manipulation and environment exploration. Given the task, AnywhereVLA parses it into simpler actions which further condition Active Environment Exploration. Exploration and navigation in larger-scale indoor environments are performed within a 3D point cloud semantic map. By leveraging a purpose-built pick-and-place dataset, AnywhereVLA exhibits robust generalization capacities.

specific manipulation and computes navigation trajectories via language-conditioned exploration algorithm, directly actuating the wheels of the mobile base and the joints of the manipulator. Our main contribution is the following:

We propose a cohesive modular framework that accepts a single language-based task instruction as input, which conditions the environment exploration and navigation modules and simultaneously drives the VLA model for manipulation task execution. Our system achieves real-time performance exceeding 10 Hz across all modules that are deployed on consumer-available edge computing units, ensuring efficient and responsive operation in dynamic settings.

## II. RELATED WORK

Large Vision-Language Models (VLMs) have emerged as a transformative paradigm. Pretraining on web-scale image-text corpora enables robust alignment between visual content and natural language semantics. When incorporated into robotic systems, these models have catalyzed a new class of policies, Vision-Language-Action (VLA) models, which extend perception-language alignment to action generation by mapping open instructions and visual context to executable behaviors [16]. VLA models unify perception, natural language, and control into end-to-end policies for robots. They leverage pretrained vision-language models and robotic trajectory datasets to ground instructions into low-level actions. MoManipVLA [13] adapts pretrained VLA models to mobile manipulators by jointly planning base and arm motions, enabling complex household operations. $\pi_0$ [11] introduces a flow matching-based action generation approach that integrates Internet-scale semantic knowledge and demonstration data for zero-shot dexterous manipulation,

while $\pi_{0.5}$ [12] extends this paradigm through co-training on heterogeneous multimodal datasets, achieving strong generalization across tasks and robots. These methods excel at instruction grounding and task generalization, but have limited spatial awareness.

Beyond this, the deployment of resource-intensive VLAs on mobile platforms requires careful optimization to ensure efficiency and maintain performance within hardware limitations. SmolVLA [17] demonstrates that a 450M-parameter VLA can achieve competitive performance compared to larger models. EdgeVLA [18] further improves efficiency by eliminating autoregressive decoding for end-effector prediction, yielding a $7\times$ speedup in inference while maintaining task accuracy. TinyVLA [19] introduces a diffusion-based policy decoder and a lightweight multimodal backbone, that match the performance of much larger VLAs while being significantly faster and more data-efficient. These advances enable real-time operation on embedded devices, though the generalization reported is primarily within manipulation domains.

A rapidly expanding body of work explores diffusion transformer-based policies for language-conditioned control. RDT-1B [20] is a billion-parameter diffusion policy pretrained on multi-robot datasets, enabling broad semantic generalization and robust bimanual skills. AC-DiT [21] introduces an adaptive coordination transformer that conditions manipulation policies in the mobility context, facilitating coupled base and arm control in mobile manipulators. However, both [20] and [21] focus exclusively on manipulation and therefore require complementary modules to provide spatial memory, target discovery, and long-horizon navigation.
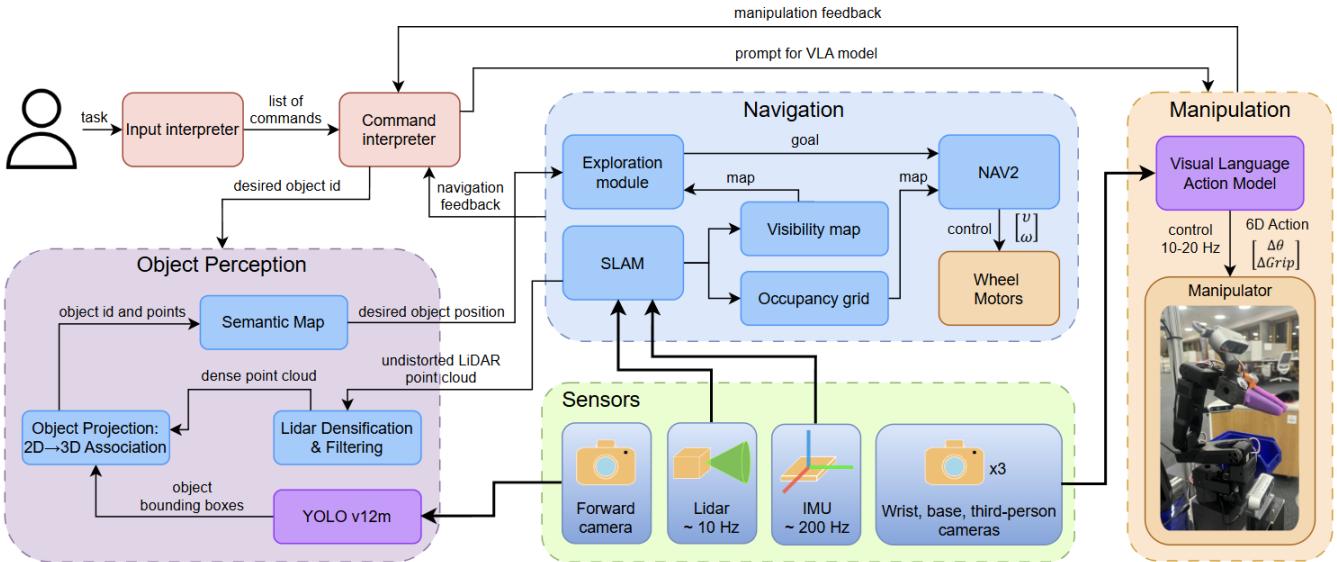
Fig. 2: AnywhereVLA architecture.

BUMBLE [14] introduces a VLM-based end-to-end framework for building-wide mobile manipulation. As a fully integrated system, it demonstrates strong performance in navigating large-scale indoor environments while manipulating a broad and previously unseen set of everyday objects. A central limitation, however, is its reliance on the known map of the whole environment and set of pre-provided landmarks. In contrast, approaches that integrate SLAM and environment exploration overcome this limitation by navigating environments autonomously without requiring prior knowledge.

Furthermore, ASC [22] frames long-horizon mobile manipulation as a sequence of modular visuomotor skills coordinated by a high-level policy. Experiments in a $185\,\mathrm{m}^2$ apartment show reliable pick-and-place with strong spatial memory, indicating effective maintenance of the task context in the rooms and for extended periods. However, complete exploration and manipulation in apartment-sized environments takes 10-15 minutes. Consequently, while ASC substantiates scalable skill coordination, its motion efficiency remains a bottleneck for deployment in real-world, human-populated spaces.

## III. ANYWHEREVLA ARCHITECTURE

AnywhereVLA framework implements an modular pipeline actuating motors of the mobile platform and manipulator's joints processing single language command and raw sensor inputs. In Fig. 2, the architecture comprises four main modules: 3D Semantic Mapping with Confidence (SM), Active Environment Exploration (AEE), Approach, and VLA Manipulation. The workflow begins with the parsing of natural language instruction, which simultaneously informs the VLA module for task-specific manipulation and conditions the AEE process. The SM module constructs a semantic 3D point cloud map by combining LiDAR-Inertial-Visual SLAM with semantic annotations from the object detection model. This map supports the AEE module, which employs frontier-based exploration conditioned on the target object class derived from the language instruction. Exploration halts once the target object is detected and localized within the semantic map. The Approach module then navigates the mobile base to the target object's location using a 2D grid map projected from a filtered LiDAR point cloud, positioning the robot for manipulation. VLA Manipulation module, leveraging a fine-tuned SmolVLA model, executes the specified task by generating actions for the manipulator's motors.

### A. 3D Semantic Mapping with Confidence (SM)

Algorithm 1 presents a 3D Semantic Object Map construction by synchronizing RGB images, undistorted LiDAR point clouds, and 2D boundbox detections by projecting LiDAR points into the camera frame to associate them with detections. To mitigate the sparsity and discontinuities of spinning LiDAR, we densified the scan by interpolating between adjacent elevation rings within each azimuth bin, followed by voxelization. For each detection, near-surface points inside its enlarged 2D bounding box are backprojected to 3D, forming object-wise point sets that are transformed to the world frame and accumulated per class over time. This produces a metric-semantic map populated by per-object 3D statistics and confidence estimates.

### B. Object Aggregation

Per class accumulated points are clustered via radius-based DBSCAN [23], outliers are robustly filtered, and each cluster is summarized by its centroid and covariance. Multi-view and data-driven cues produce an object-level confidence via a logistic mapping that fuses point density, angular coverage, inlier count, and detector scores. The resulting world-frame semantic object map exposes persistent objects with class

labels, geometry, and confidence for downstream planning and manipulation.

---

**Algorithm 1** Semantic Map Construction Pipeline

---

1: **Inputs:** image $I_t$, detections $\mathscr{D}_t$, LiDAR cloud $\mathscr{P}_t^L$, extrinsics $\mathbf{T}_{L \to C}$, TF $\mathbf{T}_{L \to W}$
2: **FOV filter:** project $\mathscr{P}_t^L$ to camera: $\mathscr{U}_t \leftarrow \Pi\big(\mathbf{T}_{L \to C}\,\mathscr{P}_t^L\big)$; keep points inside image
3: **Densification:**
4: **for** each azimuth bin $k$ **do**
5:    **for** each adjacent ring pair $(r, r+1)$ with points $\mathbf{S}_k^r, \mathbf{E}_k^{r+1}$ meeting range/gap gates **do**
6:       insert $M$ interpolants $\mathbf{P}_t = \frac{M+1-t}{M+1}\mathbf{S}_k^r + \frac{t}{M+1}\mathbf{E}_k^{r+1}$, $t=1{:}M$
7:    **end for**
8: **end for**
9: **Voxelize:** $\mathscr{Q}_t \leftarrow VoxelGrid(\text{FOV points}, v)$
10: **Projection:**
11: **for** each detection $b \in \mathscr{D}_t$ **do**
12:    enlarge box $(w, h) \leftarrow (w + \Delta_x, h + \Delta_y)$, select points in 2D box
13:    depth gate: keep points with $z \leq z_{\min}(b) + \delta$
14:    compute 3D bbox $(\mathbf{c}, \mathbf{s})$ from selected points in LiDAR frame
15:    transform points to world: $\mathscr{Q}_b^W \leftarrow \mathbf{T}_{L \to W}(t)\,\mathscr{Q}_b$
16:    accumulate $\mathscr{Q}_b^W$ into per-class database
17: **end for**
18: **Object aggregation:**
19: **for** each class $\kappa$ (periodic/on-demand) **do**
20:    aggregate points; cluster with DBSCAN$(\varepsilon, \text{minPts})$
21:    **for** each cluster $C$ **do**
22:       robust outlier filtering (MAD/$\sigma$); compute mean $\mu$ and covariance $\Sigma$
23:       estimate confidence $C$ using Eq. (2); publish markers and metadata
24:    **end for**
25: **end for**

---

### C. LiDAR densification

In Fig. 3, the sparse and discontinuous LiDAR scan pattern causes few (or no) valid returns inside a 2D detection box. In Fig. 4, we interpolate between valid ring-adjacent samples within the same azimuth bin, inserting $M$ interior points per gap (subject to range-jump and spatial-gap gates). With endpoints $\mathbf{S}$ and $\mathbf{E}$, the interpolants are

$$\mathbf{P}_t = \frac{M+1-t}{M+1}\mathbf{S} + \frac{t}{M+1}\mathbf{E}, \qquad t = 1, \dots, M . \quad (1)$$

### D. Confidence estimation

Let $\rho$ denote point density (normalized by $\rho_0$), $\Omega \in [0,1]$ denote multi-view angular coverage (union of yaw arcs over $2\pi$), $N$ the inlier count (normalized by $N_0$) and $\bar{s}$ the mean detector score. The confidence is

$$\begin{aligned} C = \sigma\Big( &w_\rho\big(1 - e^{-\rho/\rho_0}\big) + w_\Omega\,\Omega \\ &+ w_N\big(1 - e^{-N/N_0}\big) + w_S\,\bar{s} + b \Big), \end{aligned} \quad (2)$$
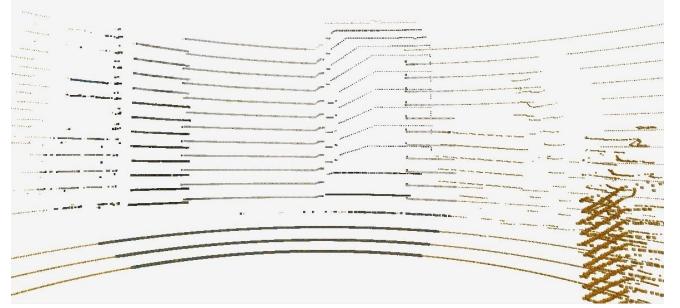


Fig. 3: Sparse point cloud from Velodyne VLP-16 LiDAR.

where $\sigma(\cdot)$ is the logistic function and $w_\rho, w_\Omega, w_N, w_S, b$ are tunable parameters.
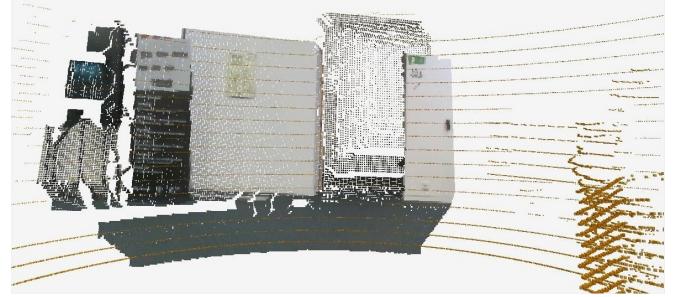


Fig. 4: Densified LiDAR point cloud after interpolation.

### E. Active Environment Exploration (AEE)

In Algorithm 2, AEE module systematically maps unknown regions to locate the target object using a frontier-based strategy [24] integrated into our modular pipeline. The occupancy grids of the SLAM backend serve as input. Frontiers are extracted through 8-connected morphological dilation and clustered. For each cluster, we compute a centroid and optimize the yaw to improve field-of-view (FoV) coverage. The resulting goal is validated with the Nav2 stack [15]. To suppress spurious goals and respect operational limits, we apply cluster filtering ($\eta_c = 20\,\text{px}$), chunking ($\eta_k = 50\,\text{px}$), non-maximum suppression ($d_{\min} = 1\,\text{m}$), an exploration radius constraint $R_e$, and FoV yaw optimization ($\alpha = 35°, R_g = 1.5\,\text{m}$). We also use a conservative frontier filter with a $5 \times 5$ kernel and a gain threshold of 50. The planner re-evaluates goals every $T_u = 4\,\text{s}$. Exploration ends upon detection of the target or upon exhaustion of valid frontiers.

### F. Approach module

The approach module computes a safe approach pose suitable for VLA manipulation from a labeled object pose on a flat support surface (e.g., a tabletop marked as an obstacle in the occupancy grid). It isolates the surface region, recovers its boundary, estimates the surface normal with principal component analysis [25], and places the robot outside the edge with a user-defined offset while orienting the base perpendicular to the surface. Nav2 [15] validates reachability and collision-free access; If the pose is infeasible, the module

**Algorithm 2** Frontier-Based Goal Selection

**Require:** Map $\mathcal{M}$, robot pose $(\mathbf{p}_r, \psi_r)$, camera params $(\alpha, R_g)$

**Ensure:** Goal pose $\mathbf{g}$ or $\perp$ (none)

1: $\mathbf{F} \leftarrow$ Dilate(free($\mathcal{M}$), $3 \times 3$) $\cap$ unknown($\mathcal{M}$) \ Dilate(occupied($\mathcal{M}$), $5 \times 5$)
2: Labels, Slices $\leftarrow$ LabelComponents($\mathbf{F}$)
3: **for** each slice $s_i$ in Slices **do**
4:     Pixels $\leftarrow$ Extract($s_i$)
5:     **if** |Pixels| $< \eta_c$ **then**
6:         **continue**
7:     **end if**
8:     Chunks $\leftarrow$ AngularPartition(Pixels, $\eta_k$)
9:     **for** each chunk $c_j$ in Chunks **do**
10:         $\mathbf{p}_{ij} \leftarrow$ Centroid($c_j$) $\oplus$ WorldTransform
11:         $\mathbf{p}_{ij} \leftarrow$ Standoff($\mathbf{p}_{ij}$, $\mathbf{p}_r$, $d_s$)
12:         **if** $\|\mathbf{p}_{ij} - \mathbf{p}_r\| > R_e$ **then**
13:             **continue**
14:         **end if**
15:         $\psi_{ij} \leftarrow$ OptimizeYaw($\mathbf{p}_{ij}$, unknown($\mathcal{M}$), $\alpha$, $R_g$)
16:         **if** Gain($\psi_{ij}$) $< 50$ **then**
17:             **continue**
18:         **end if**
19:         Add $(\mathbf{p}_{ij}, \psi_{ij})$ to Candidates
20:     **end for**
21: **end for**
22: Candidates $\leftarrow$ NMS(Candidates, $d_{\min}$); Sort by $\|\mathbf{p} - \mathbf{p}_r\|$
23: **for** each $(\mathbf{p}_i, \psi_i)$ in Candidates **do**
24:     feasible, len $\leftarrow$ ComputePath($\mathbf{p}_r$, $\mathbf{p}_i$)
25:     **if** feasible **then**
26:         $\mathbf{g} \leftarrow$ Pose($\mathbf{p}_i$, YawToQuat($\psi_i$))
27:         NavigateTo($\mathbf{g}$)
28:         **return g**
29:     **end if**
30: **end for**
31: **return** $\perp$



Fig. 5: HermesBot mobile manipulator.



(a) Third-person-view camera     (b) Wrist camera     (c) Base camera

Fig. 6: Intel Realsense D435 camera frames.

iterates over nearby edge candidates and reports failure only after exhausting valid options.

## IV. VLA MODEL FINE-TUNING

We fine-tuned the SmolVLA 450M model using an NVIDIA RTX 4090 GPU 16 GB VRAM. Fine-tuning was performed on a collected dataset comprising 50 pick-and-place episodes collected with SO-101 manipulator [26] by teleportation with a leader manipulator.

We used a batch size of 16 and a cosine decay scheduler with a learning rate of 0.0001 warmup [27]. An AdamW optimizer [28] with a weight decay of 0.01 and warmup of 100 was used. Gradients were clipped to a norm of 10.0 to mitigate the explosion.

## V. HERMESBOT MOBILE MANIPULATOR

To support the deployment of AnywhereVLA, we developed a mobile manipulator platform tailored to integrated multi-modal sensing and co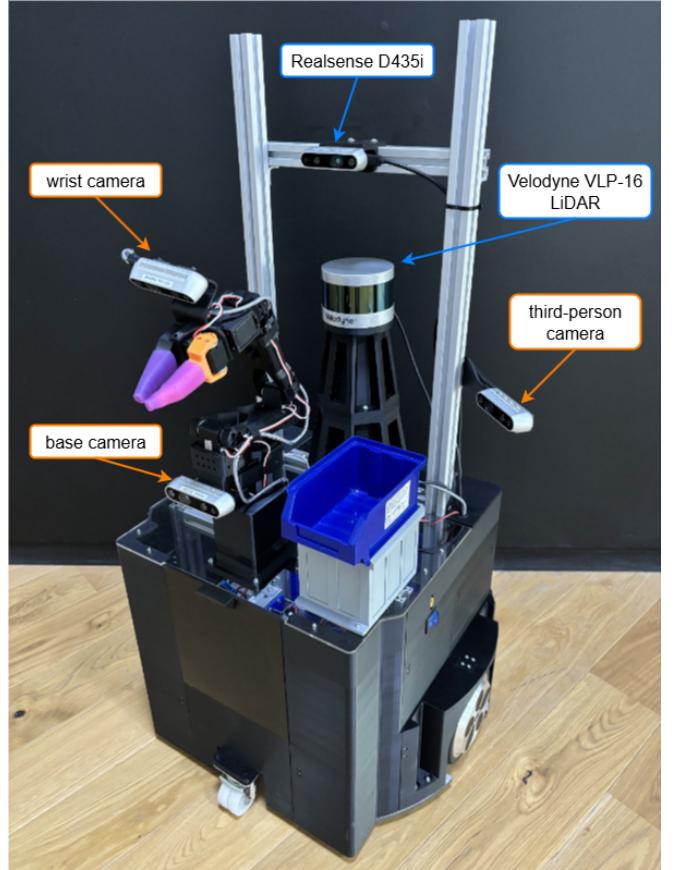mputation. In Fig. 5, The robot consists of a two-wheeled differential drive platform with the mounted SO-101 manipulator.

### A. Sensor Configuration

The sensor suite is divided into navigational and VLA subsystems. Navigation is based on a Velodyne VLP-16 LiDAR and an Intel RealSense D435i RGB-D camera, integrated with [29] for Visual-LiDAR-Inertial SLAM.

The VLA subsystem employs three Intel RealSense D435 cameras: a third-person-view camera angled to overview the manipulator and operating area (Fig. 6(a)), a wrist-mounted camera capturing the SO-101 manipulator's two-fingered gripper (Fig. 6(b)), and a base-mounted camera facing forward to monitor the manipulator's workspace (Fig. 6(c)).

## B. Computational Hardware

The computational setup comprises an NVIDIA Jetson Orin NX 16Gb for GPU-accelerated tasks and an Intel NUC Core i7 32Gb for CPU-intensive operations. The Jetson Orin processes the perception-heavy SM and VLA Manipulation modules, while the NUC handles SLAM and navigation.

Table I details the distribution of the pipeline modules, including inference frequencies and latencies measured during real-world deployment.

TABLE I: Computational throughput of AnywhereVLA modules.

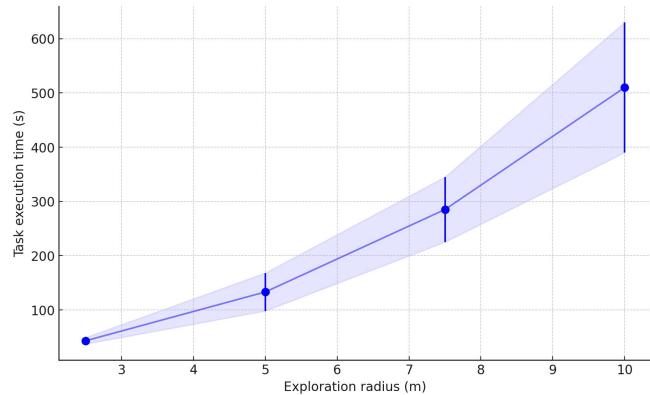| Module | Computer | Frequency (Hz) ↑ | Process time (ms) ↓ |
|---|---|---|---|
| SLAM | Intel NUC | 10 | 25 |
| Semantic Map | Jetson Orin | 15 | 45 |
| VLA | Jetson Orin | 15 | 20 |



Fig. 7: Total episode completion time.

## VI. Experimental Results

To evaluate the performance of AnywhereVLA architecture, we conducted comprehensive end-to-end experiments in diverse, unseen, indoor environments of university open space with dynamic clutter and human presence. We executed 50 pick-and-place episodes, where target objects were randomly placed within an exploration radius $R_e$. In each episode, the robot was tasked with the command: **"Pick up the `<p>object</p>` and place it in the `<p>area</p>`. And bring the `<p>object</p>` to `<p>location</p>`."** In all 50 experiments we used a blue box on HermesBot mobile manipulator as **`<p>location</p>`**. AnywhereVLA achieves an overall SR of 46% with 85% SR of VLA Manipulation module with fine-tuned SmolVLA compared to just 10% overall SR without SmolVLA fine-tuning on our purpose-built pick-and-place dataset.

Due to the modular architecture of AnywhereVLA, in Table II we highlight SR of each component individually to get insights on their influence towards the overall performance. Errors in VLA manipulation were mainly due to the bottle slipping out of the gripper and AEE failed in 25% cases being unable to find the desired object within radius $R_e$ due to cluttered tight spaces and corridors.

TABLE II: Success Rate of AnywhereVLA modules.

| Module | SR (%) ↑ |
|---|---|
| SLAM | 100 |
| Active Environment Exploration | 75 |
| Navigation | 90 |
| Object Detection | 85 |
| VLA Manipulation | 80 |

To evaluate the feasibility of real-world deployment, we measured the total episode completion time of AnywhereVLA in a set of exploration radii $R_e \in \{2.5, 5, 7.5, 10\}$ m. In Fig. 7, each run of AnywhereVLA within 5m radius, which is equivalent to an average apartment, took under 133 seconds on average. And within 10m radius our architecture managed to successfully complete tasks under 10 minutes.

## VII. Conclusion, Limitations and Future Work

In this paper, we identify a shortcoming in the existing mobile manipulation architectures: the inability to explore open and large-scale indoor environments. To address this issue, we propose an AnywhereVLA, a modular framework for language-conditioned exploration and manipulation in large-scale novel indoor environments. The system combines classical SLAM, exploration, and navigation with a lightweight VLA fine-tuned for pick-and-place operations with SO-101 by TheRobotStudio. To evaluate our system, we conducted experiments in unseen dynamic university environments. Deployed fully onboard with 3D Semantic Mapping and VLA on Jetson Orin NX and Active Environment Exploration and SLAM on Intel NUC, AnywhereVLA sustains real-time operation at $\geq 10$ Hz and achieves a 46% overall success rate completing tasks under 2.5 minutes in $80m^2$ area.

While our language-conditioned exploration pipeline demonstrates robust performance in identifying and localizing target objects within unstructured environments, AnywhereVLA exhibits a notable limitation in enforcing precise spatial-semantic constraints specified in the natural language instruction. Specifically, our pipeline is unable to successfully comprehend cases such as **"Pick up the bottle `from the table` and place it in the blue box. And bring the bottle to me."** Our architecture will explore to find the first bottle it detects, no matter where it is.

To address this limitation, future iterations of the pipeline could incorporate hierarchical semantic parsing and relational reasoning modules to disentangle object-level detection from spatial-prepositional constraints. One promising approach involves constructing a dynamic scene graph during exploration, where nodes represent detected objects and affordances, and edges encode probabilistic spatial relations derived from multimodal perception streams. Exploration policies could then be augmented with a graph-based reward signal that evaluates path efficiency and relational fidelity, penalizing deviations from the instructed configuration.. Alternatively, integrating a lightweight Vision-Language-Model

for zero-shot relational query resolution (e.g., querying "is the bottle supported by the table?") could provide an end-to-end differentiable check prior to success attribution, thereby enhancing compositional generalization without substantial computational overhead. These enhancements would align the system more closely with human-like instruction following, particularly in cluttered, multi-object scenes prevalent in real-world robotic deployment.

## REFERENCES

[1] S. Thakar, S. Srinivasan, S. Al-Hussaini, P. M. Bhatt, P. Rajendran, Y. Jung Yoon, N. Dhanaraj, R. K. Malhan, M. Schmid, V. N. Krovi, *et al.*, "A survey of wheeled mobile manipulation: A decision-making perspective," *Journal of Mechanisms and Robotics*, vol. 15, no. 2, p. 020801, 2023.

[2] J. Wu, W. Chong, R. Holmberg, A. Prasad, Y. Gao, O. Khatib, S. Song, S. Rusinkiewicz, and J. Bohg, "Tidybot++: An open-source holonomic mobile manipulator for robot learning," 2024. [Online]. Available: https://arxiv.org/abs/2412.10447

[3] Z. Fu, T. Z. Zhao, and C. Finn, "Mobile aloha: Learning bimanual mobile manipulation with low-cost whole-body teleoperation," 2024. [Online]. Available: https://arxiv.org/abs/2401.02117

[4] I. Huang, R. Cheng, S. Kim, D. Kruse, C. Matl, L. Kaul, J. Hancock, S. Harikumar, M. Tjersland, J. Borders, and D. Helmick, "Practical insights on grasp strategies for mobile manipulation in the wild," 2025. [Online]. Available: https://arxiv.org/abs/2504.12512

[5] M. Bajracharya, J. Borders, R. Cheng, D. Helmick, L. Kaul, D. Kruse, J. Leichty, J. Ma, C. Matl, F. Michel, *et al.*, "Demonstrating mobile manipulation in the wild: A metrics-driven approach," *arXiv preprint arXiv:2401.01474*, 2024.

[6] Z. He, X. Zhang, S. Jones, S. Hauert, D. Zhang, and N. F. Lepora, "Tacmms: Tactile mobile manipulators for warehouse automation," *IEEE Robotics and Automation Letters*, vol. 8, no. 8, pp. 4729–4736, 2023.

[7] C. Pu, C. Yang, J. Pu, and R. B. Fisher, "A general mobile manipulator automation framework for flexible tasks in controlled environments," *Advanced Engineering Informatics*, vol. 57, p. 102062, 2023.

[8] S. Park, X. Wang, C. C. Menassa, V. R. Kamat, and J. Y. Chai, "Natural language instructions for intuitive human interaction with robotic assistants in field construction work," *Automation in Construction*, vol. 161, p. 105345, May 2024. [Online]. Available: http://dx.doi.org/10.1016/j.autcon.2024.105345

[9] C. Zhou, H. Xu, N. Gu, Z. Wang, B. Cheng, P. Zhang, Y. Dong, M. Hayashibe, Y. Zhou, and B. He, "Language-guided long horizon manipulation with llm-based planning and visual perception," 2025. [Online]. Available: https://arxiv.org/abs/2509.02324

[10] R. Shao, W. Li, L. Zhang, R. Zhang, Z. Liu, R. Chen, and L. Nie, "Large vlm-based vision-language-action models for robotic manipulation: A survey," 2025. [Online]. Available: https://arxiv.org/abs/2508.13073

[11] K. Black, N. Brown, D. Driess, A. Esmail, M. Equi, C. Finn, N. Fusai, L. Groom, K. Hausman, B. Ichter, S. Jakubczak, T. Jones, L. Ke, S. Levine, A. Li-Bell, M. Mothukuri, S. Nair, K. Pertsch, L. X. Shi, J. Tanner, Q. Vuong, A. Walling, H. Wang, and U. Zhilinsky, "$\pi_0$: A vision-language-action flow model for general robot control," 2024. [Online]. Available: https://arxiv.org/abs/2410.24164

[12] P. Intelligence, K. Black, N. Brown, J. Darpinian, K. Dhabalia, D. Driess, A. Esmail, M. Equi, C. Finn, N. Fusai, M. Y. Galliker, D. Ghosh, L. Groom, K. Hausman, B. Ichter, S. Jakubczak, T. Jones, L. Ke, D. LeBlanc, S. Levine, A. Li-Bell, M. Mothukuri, S. Nair, K. Pertsch, A. Z. Ren, L. X. Shi, L. Smith, J. T. Springenberg, K. Stachowicz, J. Tanner, Q. Vuong, H. Walke, A. Walling, H. Wang, L. Yu, and U. Zhilinsky, "$\pi_{0.5}$: a vision-language-action model with open-world generalization," 2025. [Online]. Available: https://arxiv.org/abs/2504.16054

[13] Z. Wu, Y. Zhou, X. Xu, Z. Wang, and H. Yan, "Momanipvla: Transferring vision-language-action models for general mobile manipulation," 2025. [Online]. Available: https://arxiv.org/abs/2503.13446

[14] R. Shah, A. Yu, Y. Zhu, Y. Zhu, and R. Martín-Martín, "Bumble: Unifying reasoning and acting with vision-language models for building-wide mobile manipulation," 2024. [Online]. Available: https://arxiv.org/abs/2410.06237

[15] S. Macenski, F. Martín, R. White, and J. Ginés Clavero, "The marathon 2: A navigation system," in *2020 IEEE/RSJ Int. Conf. on Intelligent Robots and Systems (IROS)*, 2020. [Online]. Available: https://github.com/ros-planning/navigation2

[16] M. J. Kim, K. Pertsch, S. Karamcheti, T. Xiao, A. Balakrishna, S. Nair, R. Rafailov, E. Foster, G. Lam, P. Sanketi, Q. Vuong, T. Kollar, B. Burchfiel, R. Tedrake, D. Sadigh, S. Levine, P. Liang, and C. Finn, "Openvla: An open-source vision-language-action model," 2024. [Online]. Available: https://arxiv.org/abs/2406.09246

[17] M. Shukor, D. Aubakirova, F. Capuano, P. Kooijmans, S. Palma, A. Zouitine, M. Aractingi, C. Pascal, M. Russi, A. Marafioti, S. Alibert, M. Cord, T. Wolf, and R. Cadene, "Smolvla: A vision-language-action model for affordable and efficient robotics," 2025. [Online]. Available: https://arxiv.org/abs/2506.01844

[18] P. Budzianowski, W. Maa, M. Freed, J. Mo, W. Hsiao, A. Xie, T. Młoduchowski, V. Tipnis, and B. Bolte, "Edgevla: Efficient vision-language-action models," 2025. [Online]. Available: https://arxiv.org/abs/2507.14049

[19] J. Wen, Y. Zhu, J. Li, M. Zhu, K. Wu, Z. Xu, N. Liu, R. Cheng, C. Shen, Y. Peng, F. Feng, and J. Tang, "Tinyvla: Towards fast, data-efficient vision-language-action models for robotic manipulation," 2025. [Online]. Available: https://arxiv.org/abs/2409.12514

[20] S. Liu, L. Wu, B. Li, H. Tan, H. Chen, Z. Wang, K. Xu, H. Su, and J. Zhu, "Rdt-1b: a diffusion foundation model for bimanual manipulation," 2025. [Online]. Available: https://arxiv.org/abs/2410.07864

[21] S. Chen, J. Liu, S. Qian, H. Jiang, L. Li, R. Zhang, Z. Liu, C. Gu, C. Hou, P. Wang, Z. Wang, and S. Zhang, "Ac-dit: Adaptive coordination diffusion transformer for mobile manipulation," 2025. [Online]. Available: https://arxiv.org/abs/2507.01961

[22] N. Yokoyama, A. Clegg, J. Truong, E. Undersander, T.-Y. Yang, S. Arnaud, S. Ha, D. Batra, and A. Rai, "Asc: Adaptive skill coordination for robotic mobile manipulation," 2023. [Online]. Available: https://arxiv.org/abs/2304.00410

[23] M. Ester, H.-P. Kriegel, J. Sander, and X. Xu, "A density-based algorithm for discovering clusters in large spatial databases with noise," in *Knowledge Discovery and Data Mining*, 1996. [Online]. Available: https://api.semanticscholar.org/CorpusID:355163

[24] B. Yamauchi, "A frontier-based approach for autonomous exploration," *Proc. IEEE Int. Symp. Computational Intelligence in Robotics and Automation (CIRA)*, pp. 146–151, July 1997.

[25] S. Wold, K. Esbensen, and P. Geladi, "Principal component analysis," *Chemometrics and Intelligent Laboratory Systems*, vol. 2, no. 1, pp. 37–52, 1987, proceedings of the Multivariate Statistical Workshop for Geologists and Geochemists. [Online]. Available: https://www.sciencedirect.com/science/article/pii/0169743987800849

[26] R. Cadene, S. Alibert, A. Soare, Q. Gallouedec, A. Zouitine, S. Palma, P. Kooijmans, M. Aractingi, M. Shukor, D. Aubakirova, M. Russi, F. Capuano, C. Pascal, J. Choghari, J. Moss, and T. Wolf, "Lerobot: State-of-the-art machine learning for real-world robotics in pytorch," https://github.com/huggingface/lerobot, 2024.

[27] A. Gotmare, N. S. Keskar, C. Xiong, and R. Socher, "A closer look at deep learning heuristics: Learning rate restarts, warmup and distillation," 2018. [Online]. Available: https://arxiv.org/abs/1810.13243

[28] I. Loshchilov and F. Hutter, "Decoupled weight decay regularization," 2019. [Online]. Available: https://arxiv.org/abs/1711.05101

[29] C. Zheng, W. Xu, Z. Zou, T. Hua, C. Yuan, D. He, B. Zhou, Z. Liu, J. Lin, F. Zhu, Y. Ren, R. Wang, F. Meng, and F. Zhang, "Fast-livo2: Fast, direct lidar-inertial-visual odometry," 2024. [Online]. Available: https://arxiv.org/abs/2408.14035