

Bi-VLA: Vision-Language-Action Model-Based System for Bimanual Robotic Dexterous Manipulations

Koffivi Fidèle Gbagbe, Miguel Altamirano Cabrera, Ali Alabbas, Oussama Alyunes, Artem Lykov, and Dzmitry Tsetserukou

Abstract—This research introduces the Bi-VLA (Vision-Language-Action) model, a novel system designed for bimanual robotic dexterous manipulation that seamlessly integrates vision for scene understanding, language comprehension for translating human instructions into executable code, and physical action generation. We evaluated the system’s functionality through a series of household tasks, including the preparation of a desired salad upon human request. Bi-VLA demonstrates the ability to interpret complex human instructions, perceive and understand the visual context of ingredients, and execute precise bimanual actions to prepare the requested salad. We assessed the system’s performance in terms of accuracy, efficiency, and adaptability to different salad recipes and human preferences through a series of experiments. Our results show a 100% success rate in generating the correct executable code by the Language Module, a 96.06% success rate in detecting specific ingredients by the Vision Module, and an overall success rate of 83.4% in correctly executing user-requested tasks.

Keywords— *Vision-Language-Action Model, Bimanual Robotic Manipulation, Human-Robot Interaction, Generative AI.*

I. INTRODUCTION

Recent advancements in language models have significantly impacted Human-Robot Interaction (HRI), enabling robots to engage more naturally and effectively with their human counterparts [1]. In particular, to enhance performance efficiency and tackle complex tasks, coordinated dual-arm robotic systems are increasingly employed. These systems offer greater reliability and facilitate the execution of tasks that are challenging for a single manipulator [2]. Furthermore, the coordination of two manipulator systems has been extensively studied in the literature, yielding satisfactory performance outcomes [3], [4]. However, traditional planning-based methods, which focus on motion planning through predefined trajectories for coordinating two robotic arms [5], [6], are often suboptimal for dynamic and complex tasks. To achieve human-level bimanual manipulation, alternative strategies such as reinforcement learning (RL) [7] and learning from demonstration (LfD) [8] have been proposed. Nevertheless, it is crucial to acknowledge the inherent challenges associated with these approaches, particularly regarding the time required for learning and data collection.

Moreover, the integration of language models in HRI not only facilitates natural communication but also extends their utility beyond basic interactions. This includes tasks

The authors are with the Intelligent Space Robotics Laboratory, Center for Digital Engineering, Skolkovo Institute of Science and Technology, Moscow, Russia {Koffivi.Gbagbe, M.Altamirano, Ali.Alabbas, Oussama.Alyunes, Artem.Lykov, D.Tsetserukou}@skoltech.ru

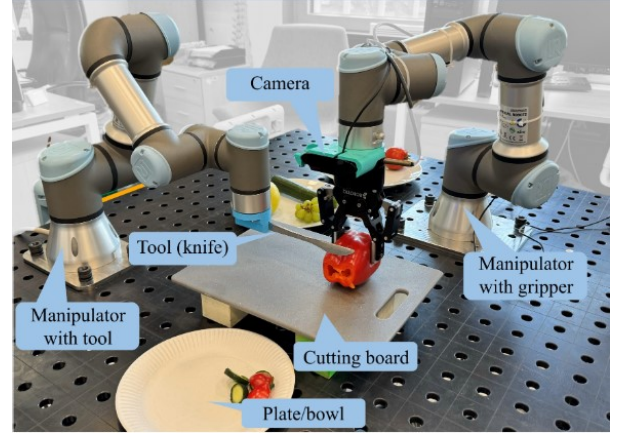


Fig. 1: System Overview of Bi-VLA.

such as task planning and generating intermediate reasoning steps, commonly referred to as a chain of thought. However, despite their potential, the application of language models to synthesize the bimanual skills of robots has not received significant attention. This is particularly noteworthy given that numerous tasks humans perform routinely require the coordinated use of both hands, highlighting the essential role of bimanual capabilities in executing intricate operations that are beyond the reach of single-arm manipulators. Building on this premise, the integration of vision, language, and action within a unified model presents a novel approach to enhancing robotic bimanual operation. In this research, we aim to develop a vision-language-action model, dubbed Bi-VLA, specifically designed for bimanual robotic operations. Additionally, we will explore its potential impact on coordinating the actions of dual-arm robotic systems in response to human instructions. The proposed scenario is a cooking robot, where a robot is in charge of picking and placing the required ingredients, and the second cutting them.

II. RELATED WORK

A. From language instructions to robotic actions

Translating user language input into robotic actions has seen considerable progress through the integration of Large Language Models (LLMs) and various algorithmic frameworks. CLIPort [9] leverages visual and language models for robotic manipulation, focusing on the spatial understanding necessary for accomplishing the tasks. Similarly, ProgPrompt [10] introduced a programmatic LLM prompt structure that enables plan generation for robotic tasks scoring potential

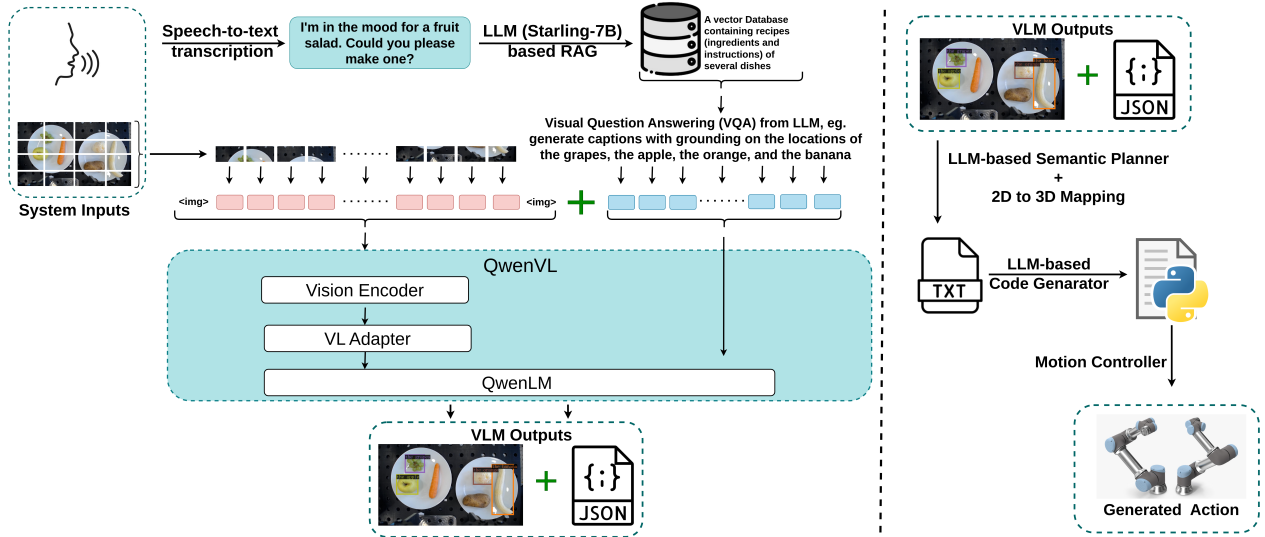


Fig. 2: Bi-VLA Architecture.

actions without the need for explicit domain knowledge. The “Code as Policies” [11] paradigm further underscores the importance of programmatic structures in embodying LLM-generated instructions for control, which is in close alignment with RT-1’s objective of scalable, real-world robotic control [12]. Drawing on pre-trained vision-language models, open-world object manipulation [13] is evolving to adapt robotic actions to diverse and changing environments. Instruct2Act [14] steps beyond basic instructions by aligning multimodal inputs with actionable robotic sequences, aided by the breadth of responses offered by LLMs. “Language to Rewards” [15] introduced a method of reward function generation from natural language inputs to calibrate and optimize robotic skills, with the goal of building cooperative agents by harnessing modularity through LLM architectures. Continuing this thread, RT-2 [16] focuses on translating visual perceptions and user language into actionable instructions, demonstrating the critical intersection between understanding and physical task execution. Lastly, PaLM-E[17] highlights the convergence of multimodal understanding within LLMs, emphasizing the embodiment of language models in tangible robotic actions, while LLM-MARS[18] incorporates LLMs into the multi-agent robotics system for behavior trees generation. Collectively, these advancements contribute to an overarching framework that not only enhances robots’ ability to understand and act on complex user language inputs but also redefines interactions between humans and robots. In this work, we use LLMs as a strategic semantic planner, orchestrating the collaborative efforts of robotic arms to successfully complete the given task.

B. Bimanual robotic manipulation

The exploration of bimanual object manipulation spans a multidisciplinary field, incorporating insights from neuroscience, robotics, and computer science. To establish a foundation, researchers have proposed various taxonomies aimed at emulating human-like coordination in complex

tasks. Notably, the integration of spatial and temporal constraints is essential for achieving dexterity [19]. Moreover, behavioral control paradigms suggest that dual-arm coordination often surpasses single-arm methods, highlighting the necessity for behavior-based control systems to facilitate effective bimanual interaction, as demonstrated by studies utilizing deep learning for imitation learning that illustrate effective teaching methods for robots through the observation and replication of human actions [20]. Additionally, low-cost teleoperation solutions, such as Mobile ALOHA [21], have opened avenues for learning complex manipulation behaviors. Finally, progress in this area is complemented by advancements in sim-to-real deep reinforcement learning, which refines bimanual tactile manipulation abilities in simulated environments before transferring them to real-world applications [22]. In this work, we focus on a less data-intensive approach that eliminates the need for training or additional fine-tuning. This method involves the language module invoking specific predefined action APIs, determining their calling sequence, and allocating tasks between the two robotic arms to complete the user’s assigned task.

III. GROUNDING LANGUAGE INSTRUCTIONS INTO BIMANUAL ROBOT ACTIONS

A. General system overview

The system takes as input the user’s task request, which is transcribed into a text message, along with an image from the camera (see Fig. 2 left). The task request serves as a query for the Retrieval-Augmented Generation (RAG) system, and its output initiates a visual question-answering (VQA) dialogue between the Language Module and the Vision Language Module. The input image, on the other hand, is split into patches, tokenized, and combined with the VQA question tokens before being sent to the Vision Language Module. The main outputs include an image with captions and a JSON file containing the 2D positions of ingredients, a list of all ingredients, and missing ingredients, etc. Subsequently,

the Language Module generates a set of actions, which are performed by the robots (see Fig. 2 right). In Fig. 3 an example of a set of actions can be observed when all the required ingredients are available. In Fig. 4 the details of a single action are shown.

B. Large Language Model Module

In this work, we utilized **Starling-LM-7B-alpha** [23], an open-source LLM developed by UC Berkeley’s NEST Lab using Reinforcement Learning from AI Feedback (RLAIF), as the backbone for the language module, which outperforms **GPT-3.5** [24] in several benchmarks, including an 8.09 out of 10 in MT Bench.

1) *Semantic planning*: Semantic planning involves the use of natural language understanding and generation techniques to interpret human instructions and generate executable plans. In the context of our work, the semantic planner is designed to generate a plan for coordinating the movements of the two robot arm manipulators equipped with different and appropriate tools in order to perform a given action. The planner takes the action description as input and produces a plan specifying the sequence of motions that need to be executed by the robot arms (see Fig. 4). The planner utilizes a template that provides a general structure for the plan. It starts with the “[start of plan]” tag and ends with the “[end of plan]” tag, ensuring that the plan is clearly delineated. The template also includes specific instructions and placeholders that need to be filled based on the task requirements by the vision-language model (VLM) (see section III-C.1), as well as a set of rules to follow.

2) *Translating semantic plan into APIs function calls with the Code Generator*: The plan generated by the semantic planner is translated into a sequence of API call functions. The goal is to convert the textual plan into executable Python code that adheres to a predefined set of functions provided by the API. We use a prompt that outlines the functions available for use, the structure of the plan, and the rules to follow. To illustrate the process, we provide an example plan and its corresponding Python code (Fig. 4), helping to guide the model’s response as in the one-shot prompting technique. This methodology facilitates the seamless translation of high-level task plans into executable Python code, promoting efficiency and clarity.

3) *From API function call to Action Generation with The Motion Control APIs*: To facilitate action generation, we have implemented a comprehensive set of API functions to control the robot’s actions. These functions are designed to perform a variety of tasks, including:

- **open_gripper(manipulator_type)**: to open the gripper.
- **move_to_object(manipulator_type, object_name)**: to move the manipulator to the position of an object.
- **grasp(manipulator_type, object_name)**: to grasp an object.
- **cut(manipulator_type, object_name)**: for cutting actions.
- **pour(manipulator_type, object_name)**: for pouring actions.

- **put(manipulator_type, object_name)**: for pick-and-place actions.
- **toss(manipulator_type, object_name)**: for mixing objects or vegetables in a bowl.
- **cut_and_put_in(manipulator_type, object_name)**: which combines the **cut** and **put** actions.

The vision language model uses the function **get_list_of_objects()** to retrieve the names of all objects visible on the table within the camera frame, as well as **get_bounding_boxes()** which generates the bounding boxes of the desired list of objects. Ensuring explicit naming of API functions is crucial to preventing confusion with the language model. This comprehensive set of API functions provides the necessary control and versatility for robotic operations, enabling a wide range of complex tasks to be performed with precision and efficiency.

C. Vision Language Model Module

1) *Integration of the Qwen Vision-Language Model*: The QWEN-VL series are large-scale vision-language models (VLMs) that understand both text and images [25]. Built on the QWEN-LM foundation, they feature visual capabilities through a specialized architecture and training pipeline, achieving state-of-the-art performance across various benchmarks, including image captioning, visual question answering, and visual grounding. In this work, we use the model to verify the availability of required items and provide the 2D-pixel coordinates of these items.

2) *Mapping Image Pixel to 3D Coordinates*: We assumed that the position of an object is equal to the position of its center. The relationship between the position of an object in the world coordinate system and pixel coordinates is given by the equation:

$$z_c {}^pP = {}^pT_c {}^c\hat{T}_w {}^wP, \quad (1)$$

where z_c is the position of the object on the z axis of the camera coordinates; ${}^pP = [x_p \ y_p \ 1]^T$ is the pixel coordinate of the object; pT_c is the intrinsic matrix of the camera, a 3×3 matrix, that transforms position vectors from the camera coordinates to the pixel coordinates, and we get it when we calibrate the camera; ${}^c\hat{T}_w$ is the 3×4 matrix that transforms position vectors from world coordinates to camera coordinates, and it is known since the camera is fixed on the robot tool and the transformation between the tool and the world coordinates is known all the time, and ${}^wP = [x_w \ y_w \ z_w \ 1]^T$ is the position of the object in the world coordinates.

We rewrite ${}^c\hat{T}_w$ as:

$${}^c\hat{T}_w = \left[\begin{array}{c|c} {}^cR_{3 \times 3} & {}^c_t t_{3 \times 1} \end{array} \right], \quad (2)$$

where cR is the rotation matrix from the world coordinates to the camera coordinates and c_t is the translation vector. We also define

$${}^cT_w = \left[\begin{array}{c|c} {}^c\hat{T}_w \\ \hline 0 & 0 & 0 & 1 \end{array} \right] = \left[\begin{array}{c|c} {}^cR_{3 \times 3} & {}^c_t t_{3 \times 1} \\ \hline 0_{1 \times 3} & 1 \end{array} \right]. \quad (3)$$

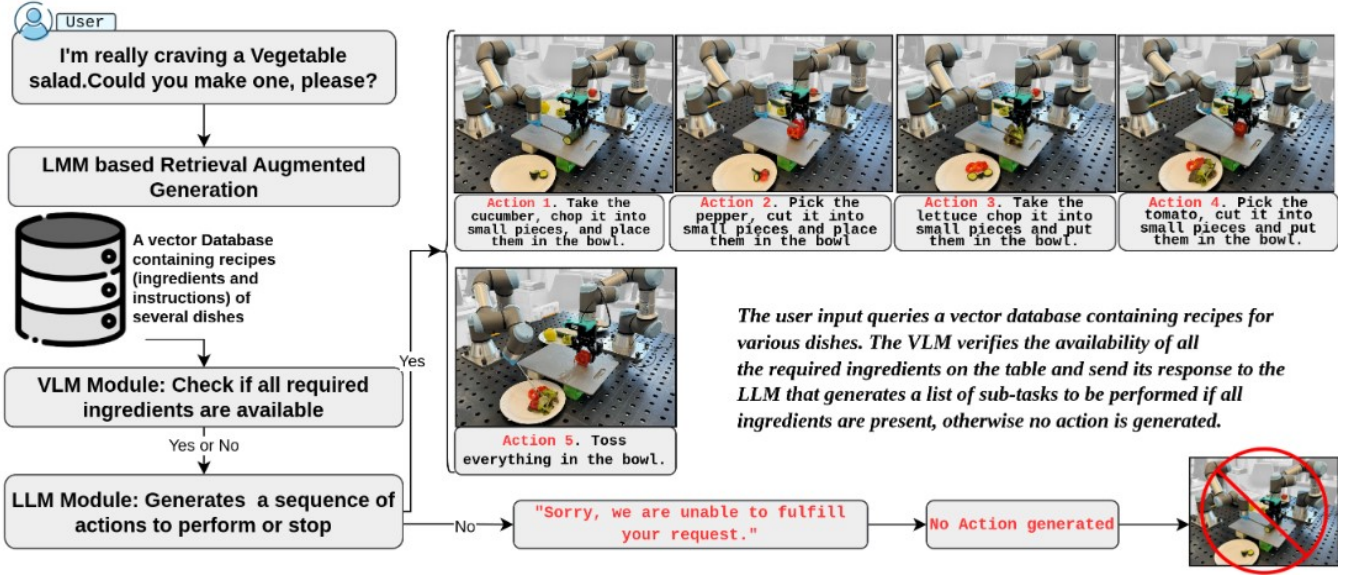


Fig. 3: Cooking Experiment Workflow.

The pixel position that we get from the camera is distorted with radial and tangential errors [26]. In order to obtain accurate transformation from pixel to real-world coordinates, these distortion errors have to be taken into account. When the camera was calibrated, we obtained the intrinsic matrix p_cT and the distortion coefficients $[k_1, k_2, p_1, p_2, k_3]$. To correct the distortion error, we used the Brown–Conrady model that follows the following equations [27]:

$$\begin{aligned} x_u &= x_d(1 + k_1r^2 + k_2r^4 + k_3r^6) + \\ &\quad + 2p_1x_dy_d + p_2(r^2 + 2x_d^2), \\ y_u &= y_d(1 + k_1r^2 + k_2r^4 + k_3r^6) + \\ &\quad + p_1(r^2 + 2y_d^2) + 2p_2x_dy_d, \end{aligned} \quad (4)$$

where x_u and y_u are the undistorted pixel coordinates, x_d and y_d are the distorted pixel coordinates, $[k_1, k_2, p_1, p_2, k_3]$ are the distortion coefficients and $r = \sqrt{x_d^2 + y_d^2}$ is the Euclidean distance of the distorted point.

To transfer the position of the objects from Pixel to 3D coordinates, we followed the following steps:

- 1) Get the distorted pixel position of the center of the object $[x_d \ y_d]^T$.
- 2) Calculate the undistorted pixel coordinates $[x_u \ y_u]^T$ from the distorted ones using Eq. 4.
- 3) Calculate the normalized coordinates in the camera coordinates system depending on Eq. 1.

$$\begin{bmatrix} x_c \\ y_c \\ z_c \\ 1 \end{bmatrix} = {}^p_cT^{-1} \begin{bmatrix} x_u \\ y_u \\ 1 \end{bmatrix},$$

where x_c, y_c, z_c are the position of the object center in the camera coordinates system.

- 4) Calculate the real position in the camera coordinates, $[x_c \ y_c \ z_c]^T$, from the normalized vector by multiplying it by z_c . In our application all objects are placed on

the table in the same z –plane relative to the camera coordinates system and the camera is operating in a known altitude from the table and its z axis is perpendicular to the table.

- 5) Calculate the position of the object in the world coordinates $[x_w \ y_w \ z_w]$ from the position in the camera coordinates.

$$\begin{bmatrix} x_w \\ y_w \\ z_w \\ 1 \end{bmatrix} = {}^c_wT^{-1} \begin{bmatrix} x_c \\ y_c \\ z_c \\ 1 \end{bmatrix}$$

These steps allow us to get the world coordinates of the objects, which are then sent to the robot. After receiving the position, the robot moves its Tool Center Point (TCP) above the object and then gets it down to catch the object.

IV. EXPERIMENTAL EVALUATION

A. Experimental Setup

Two UR3 collaborative robots from Universal Robots were mounted on a Siegmund welding table (see Fig. 1). A custom-designed and 3D-printed camera holder, made of PLA, was attached to the first manipulator, which also held a Logitech HD 1080p webcam and a two-finger Robotiq gripper. A knife holder was similarly designed and installed in the second manipulator's end effector. Ingredients for various salads were positioned near the first manipulator where the camera was able to capture images to identify the necessary items. A cutting board was placed between the two manipulators for cutting operations, with a plate or bowl in front to place the prepared ingredients. All processes are supervised by a PC (NVIDIA GeForce RTX 4090 24Gb, SSD 2Tb, RAM 64Gb 5600MHz, Intel Core i9, running Ubuntu 22.04).

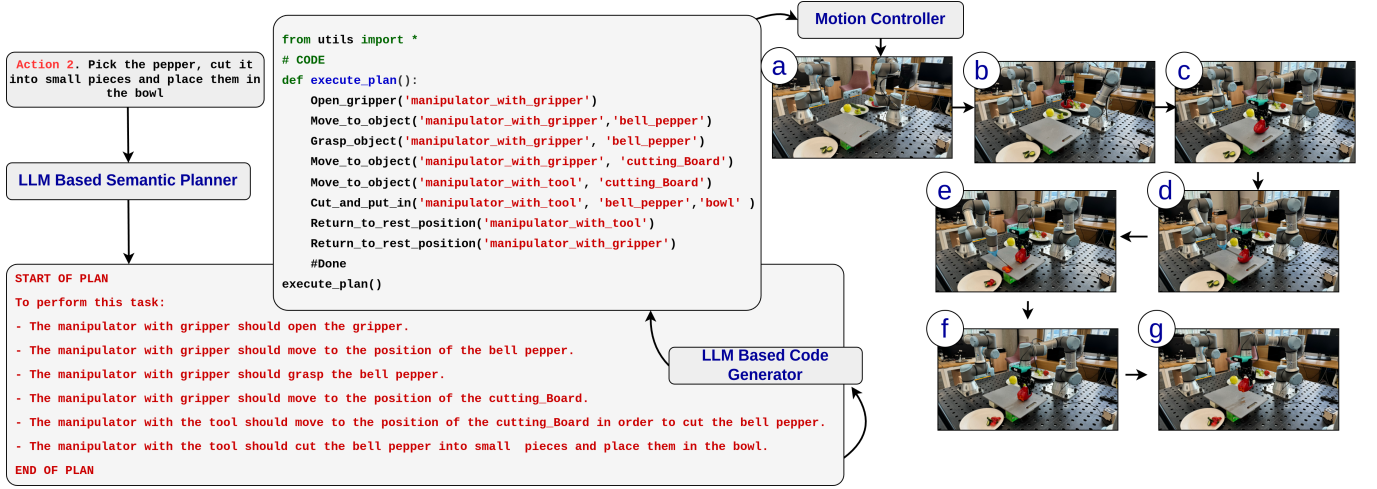


Fig. 4: Example illustration of Action 2 in Fig. 3. The LLM-based semantic planner receives the action to perform and generates a detailed plan outlining the movements of two robot manipulators. The generated plan is translated into a set of motion API call functions through the LLM-based Code Generator. The execution of the generated code triggers the movement of the two robot manipulators through the Motion Controller. (a)-(b) the manipulator with the gripper moves to grasp the pepper and bring it to the cutting board; (c)-(d) the manipulator with the knife moves to the cutting board, cuts the pepper (e), and places it in the bowl; (f)-(g) the two manipulators return to their initial position.

B. Experiments Workflow

As depicted in Fig. 3, we design several custom cooking recipes stored in a vector database, which serves as a set of tasks for the robot manipulators to perform collaboratively. The recipes clearly state the required ingredients as well as the instructions to follow for successful execution. Among these sets of tasks, we evaluate the performance of our system on three selected tasks, including making a vegetable salad, a Russian salad, and a fruit salad. The three tasks are different in nature in the sense that they require several types of actions, such as picking, placing, cutting, pouring, mixing, etc. When a user submits a request, it serves as a query processed through the LLM-Based Retrieval-Augmented Generation (RAG) technique [28], which enhances response quality and factuality by augmenting the language model’s outputs with relevant information from external knowledge sources, in this case, a vector database. After the RAG model produces a response, the VLM module verifies the availability of all the ingredients mentioned in the recipe. Based on the VLM’s response, the LLM module will then generate a list of actions for the two robots to perform if the ingredients are available. Alternatively, if the ingredients are not available or the request is ambiguous, no action will be performed.

V. RESULTS AND EVALUATION

We assessed the performance of our system’s components both individually and collectively, enabling us to identify their strengths and weaknesses and understand their interactions in contributing to overall performance. We define a task as completed successfully if the generated code by the Code Generator is correct before execution by the Motion Controller and if the actions are successfully executed.

- 1) Firstly, we evaluated the performance of the Language Module by generating 100 different requests with GPT-3.5 for the system to prepare three types of salad, e.g. “I’m really craving a vegetable salad. Could you make one, please?”. Specifically, one-third of the requests were for making vegetable salads, another third were for a Russian salad, and the final third was for a fruit salad. Our proposed Language Module successfully generated the correct executable code for all the requests, achieving a success rate of 100%.
- 2) In the second part of the experiment, we evaluated the vision module’s performance using 100 pictures of different ingredients, with one-third allocated to each salad type. For each salad type, half the pictures included all the required ingredients (Scenario 1), while the other half were missing some (Scenario 2). We evaluated if the model returned the correct list of all the ingredients in the input image and whether each mentioned ingredient from the salad recipe was detected correctly; that is, the captions on the ingredients and the 2D bounding boxes were accurate in the output image.
 - In Scenario 1, where all ingredients were correct, we achieved a 96.06% success rate in detecting ingredients and an 83.4% success rate in generating correct actions.
 - In Scenario 2, the vision module misidentified present ingredients as missing ones, leading to a detection performance of 86.53% and a 71.22% success rate in generating correct actions.
- 3) For the Action Generation Module, we ensured reliable translation from 2D pixel to 3D world coordinates as described in Section III-C.2, enabling the robot’s grip-

per to successfully grasp objects of varying orientations and perform cutting operations accurately. Therefore, we assume that the success rate of the module is equivalent to the success rate of generating correct 2D bounding box coordinates.

The table below summarizes the main experiment results:

TABLE I: Experiment Results

Experiments	Success rate (%)
LLM Module (all scenario)	100%
VLM Module (scenario 1)	96.06%
VLM Module (scenario 2)	86.53%
Action Generation Module (scenario 1)	83.4%
Action Generation Module (scenario 2)	71.22%

VI. CONCLUSION

In this work, we explore a new approach to dexterous bimanual manipulation by integrating vision, language understanding, and robotic action generation into a unified system named Bi-VLA. Our experiments demonstrate the system’s capability to accurately interpret human instructions, perceive the visual context of the scene, and execute the required actions. The high success rates of the language and vision modules emphasize their crucial roles in overall system performance. Our findings highlight the need for ongoing advancements in visual understanding, as the accuracy of the vision module significantly affects task completion. Future efforts will focus on enhancing the robustness and versatility of the vision module to better handle a variety of visual variations and uncertainties, ultimately making robotic assistants more reliable and adaptable.

ACKNOWLEDGEMENTS

Research reported in this publication was financially supported by the RSF grant No. 24-41-02039.

REFERENCES

- [1] C. Y. Kim, C. P. Lee, and B. Mutlu, “Understanding large-language model (llm)-powered human-robot interaction,” in *Proc. ACM/IEEE Int. Conf. on Human-Robot Interaction (HRI)*, 2024, pp. 371–380.
- [2] C. D. Vo, D. A. Dang, and P. H. Le, “Development of multi-robotic arm system for sorting system using computer vision,” *Journal of Robotics and Control (JRC)*, vol. 3, no. 5, pp. 690–698, 2022.
- [3] A. Edsinger and C. C. Kemp, “Two arms are better than one: A behavior based control system for assistive bimanual manipulation,” in *Recent Progress in Robotics: Viable Robotic Service to Human: An Edition of the Selected Papers from the 13th International Conference on Advanced Robotics*. Springer, 2007, pp. 345–355.
- [4] D. Rakita, B. Mutlu, M. Gleicher, and L. M. Hiatt, “Shared control-based bimanual robot manipulation,” *Science Robotics*, vol. 4, no. 30, 2019.
- [5] C. Bersch, B. Pitzer, and S. Kammel, “Bimanual robotic cloth manipulation for laundry folding,” in *Proc. IEEE/RSJ Int. Conf. on Intelligent Robots and Systems*, 2011, pp. 1413–1419.
- [6] S. S. Mirrazavi Salehian, N. B. Figueroa Fernandez, and A. Billard, “Dynamical system-based motion planning for multi-arm systems: Reaching for moving objects,” in *Proc. of the 26th Int. Joint Conf. on Artificial Intelligence (IJCAI)*, 2017, pp. 4914–4918.
- [7] Y. Chen, T. Wu, S. Wang, X. Feng, J. Jiang, Z. Lu, S. McAleer, H. Dong, S.-C. Zhu, and Y. Yang, “Towards human-level bimanual dexterous manipulation with reinforcement learning,” *Advances in Neural Information Processing Systems*, vol. 35, pp. 5150–5163, 2022.

- [8] G. Franzese, L. de Souza Rosa, T. Verburg, L. Peternel, and J. Kober, “Interactive imitation learning of bimanual movement primitives,” *IEEE/ASME Transactions on Mechatronics*, 2023.
- [9] M. Shridhar, L. Manuelli, and D. Fox, “Cliport: What and where pathways for robotic manipulation,” in *Proc. of the 5th Conf. on Robot Learning (CoRL)*. PMLR, 2022, pp. 894–906.
- [10] I. Singh, V. Blukis, A. Mousavian, A. Goyal, D. Xu, J. Tremblay, D. Fox, J. Thomason, and A. Garg, “Progprompt: Generating situated robot task plans using large language models,” in *Proc. IEEE Int. Conf. on Robotics and Automation (ICRA)*. IEEE, 2023, pp. 11 523–11 530.
- [11] J. Liang, W. Huang, F. Xia, P. Xu, K. Hausman, B. Ichter, P. Florence, and A. Zeng, “Code as policies: Language model programs for embodied control,” in *Proc. IEEE Int. Conf. on Robotics and Automation (ICRA)*. IEEE, 2023, pp. 9493–9500.
- [12] A. Brohan, N. Brown, J. Carbajal, Y. Chebotar, J. Dabis, C. Finn, K. Gopalakrishnan, K. Hausman, A. Herzog, J. Hsu *et al.*, “Rt-1: Robotics transformer for real-world control at scale,” *arXiv preprint arXiv:2212.06817*, 2022.
- [13] A. Stone, T. Xiao, Y. Lu, K. Gopalakrishnan, K.-H. Lee, Q. Vuong, P. Wohlhart, S. Kirmani, B. Zitkovich, F. Xia *et al.*, “Open-world object manipulation using pre-trained vision-language models,” *arXiv preprint arXiv:2303.00905*, 2023.
- [14] S. Huang, Z. Jiang, H. Dong, Y. Qiao, P. Gao, and H. Li, “Instruct2act: Mapping multi-modality instructions to robotic actions with large language model,” *arXiv preprint arXiv:2305.11176*, 2023.
- [15] W. Yu, N. Gileadi, C. Fu, S. Kirmani, K.-H. Lee, M. G. Arenas, H.-T. L. Chiang, T. Erez, L. Hasenclever, J. Humpik *et al.*, “Language to rewards for robotic skill synthesis,” *arXiv preprint arXiv:2306.08647*, 2023.
- [16] A. Brohan, N. Brown, J. Carbajal, Y. Chebotar, X. Chen, K. Choro-manski, T. Ding, D. Driess, A. Dubey, C. Finn, P. Florence *et al.*, “Rt-2: Vision-language-action models transfer web knowledge to robotic control,” in *arXiv preprint arXiv:2307.15818*, 2023.
- [17] D. Driess, F. Xia, M. S. Sajjadi, C. Lynch, A. Chowdhery, B. Ichter, A. Wahid, J. Tompson, Q. Vuong, T. Yu *et al.*, “Palm-e: An embodied multimodal language model,” *arXiv preprint arXiv:2303.03378*, 2023.
- [18] A. Lykov, M. Dronova, N. Naglov, M. Litvinov, S. Satskevich, A. Bazhenov, V. Berman, A. Shcherbak, and D. Tssetserukou, “Llm-mars: Large language model for behavior tree generation and nlp-enhanced dialogue in multi-agent robot systems,” *arXiv preprint arXiv:2312.09348*, 2023.
- [19] F. Krebs and T. Asfour, “A bimanual manipulation taxonomy,” *IEEE Robotics and Automation Letters*, vol. 7, no. 4, pp. 11 031–11 038, 2022.
- [20] F. Xie, A. Chowdhury, M. De Paolis Kaluza, L. Zhao, L. Wong, and R. Yu, “Deep imitation learning for bimanual robotic manipulation,” *Advances in neural information processing systems*, vol. 33, pp. 2327–2337, 2020.
- [21] Z. Fu, T. Z. Zhao, and C. Finn, “Mobile aloha: Learning bimanual mobile manipulation with low-cost whole-body teleoperation,” *arXiv preprint arXiv:2401.02117*, 2024.
- [22] Y. Lin, A. Church, M. Yang, H. Li, J. Lloyd, D. Zhang, and N. F. Lepora, “Bi-touch: Bimanual tactile manipulation with sim-to-real deep reinforcement learning,” *IEEE Robotics and Automation Letters*, vol. 8, no. 9, pp. 5472–5479, 2023.
- [23] B. Zhu, E. Frick, T. Wu, H. Zhu, and J. Jiao, “Starling-7b: Improving llm helpfulness & harmlessness with rlai,” November 2023.
- [24] G. B. et al., “Introducing chatgpt,” 2022. [Online]. Available: <https://openai.com/blog/chatgpt>
- [25] J. Bai, S. Bai, S. Yang, S. Wang, S. Tan, P. Wang, J. Lin, C. Zhou, and J. Zhou, “Qwen-vl: A frontier large vision-language model with versatile abilities,” *arXiv preprint arXiv:2308.12966*, 2023.
- [26] J. P. De Villiers, F. W. Leuschner, and R. Geldenhuys, “Centi-pixel accurate real-time inverse distortion correction,” in *Proc. Int. Symp. on Optomechatronic Technologies*, vol. 7266. SPIE, 2008, pp. 320–327.
- [27] A. E. Conrady, “Decentred lens-systems,” *Monthly notices of the royal astronomical society*, vol. 79, no. 5, pp. 384–390, 1919.
- [28] P. Lewis, E. Perez, A. Piktus, F. Petroni, V. Karpukhin, N. Goyal, H. Küttler, M. Lewis, W. tau Yih, T. Rocktäschel, S. Riedel, and D. Kiela, “Retrieval-augmented generation for knowledge-intensive nlp tasks,” in *Proc. Int. Conf. on Neural Information Processing Systems (NIPS)*, 2021, p. 16.