# MOSAIC: Modular Foundation Models for Assistive and Interactive Cooking

**Huaxiaoyue Wang**\*, **Kushal Kedia**\*, **Juntao Ren**\*,
**Rahma Abdullah, Atiksh Bhardwaj, Angela Chao, Kelly Y Chen,**
**Nathaniel Chin, Prithwish Dan, Xinyi Fan, Gonzalo Gonzalez-Pumariega,**
**Aditya Kompella, Maximus Adrian Pace, Yash Sharma, Xiangwan Sun, Neha Sunkara,**
**Sanjiban Choudhury**
Cornell University

**Abstract:** We present MOSAIC, a modular architecture for coordinating multiple robots to *(a)* interact with users using natural language and *(b)* manipulate an open vocabulary of everyday objects. MOSAIC employs modularity at several levels: it leverages multiple large-scale pre-trained models for high-level tasks like language and image recognition, while using streamlined modules designed for low-level task-specific control. This decomposition allows us to reap the complementary benefits of foundation models as well as precise, more specialized models. Pieced together, our system is able to scale to complex tasks that involve coordinating multiple robots and humans. First, we unit-test individual modules with 180 episodes of visuomotor picking, 60 episodes of human motion forecasting, and 46 online user evaluations of the task planner. We then extensively evaluate MOSAIC with 60 end-to-end trials. We discuss crucial design decisions, limitations of the current system, and open challenges in this domain. The project's website is at https://portal-cornell.github.io/MOSAIC/

**Keywords:** Robot Learning, Foundation Models, Human-Robot Interaction

## 1 Introduction

Collaborative tasks in household environments present significant challenges for robots. Consider the scenario in Figure 1, where a human user collaborates with two robots to prepare a meal. We'd like for communication with the system should feel natural to the user. Furthermore, each robot should be able to complete auxiliary tasks across a wide range of objects while fluidly collaborate with the humans. Prior systems in this domain [1, 2, 3, 4, 5] have demonstrated impressive capabilities. However, they either they operate in isolation and lack meaningful collaboration with humans, or employ highly scripted behavior. In this paper, we aim to overcome both of these limitations by designing a system that fluidly collaborates with humans and performs a wide range of tasks.

While a single end-to-end model works well for tasks like language understanding where large amounts of data are available, such an approach is difficult for collaborative robots, where less data is available and precise control is important. Our key insight is that ***by modularizing our architecture, we can segment out parts of the framework that require broad generalization, such as language and image recognition, from the portions that require task-specific control. Furthermore, we can triage end-to-end system failures to a specific component, enabling efficient system improvement.***

We operationalize this insight to create MOSAIC (**Mo**dular **S**ystem for **A**ssistive and **I**nteractive **C**ooking): an architecture that applies modularity at *multiple distinct levels* to significantly improve the overall system's performance. Each module has a well-scoped task which is simpler to complete

---

\* Denotes equal contribution.
Correspondence to: Huaxiaoyue Wang, yukiwang@cs.cornell.edu

Figure 1: **MOSAIC cooking in the kitchen.** MOSAIC interacts with a user via natural language and controls a tabletop manipulator (R1) and a mobile manipulator (R2) to prepare vegetable soup with the user.

and results in fewer overall mistakes. While the principle of modularity has been central to developing robust real-world robotic systems, these modules have historically been robot/task specific. In contrast, our architecture integrates general-purpose pre-trained models to solve robotic tasks. This design choice enables us to build a system that is flexible, interpretable, and scalable.

Our contributions can be organized into three groups:

**1. We architect MOSAIC: a full-stack Modular Cooking Assistant.** We present a novel framework for home robots that integrates multiple large-scale pre-trained models. In particular, we use large language models for interactive task planning, vision language models for visuomotor skills, and motion forecasting models for predicting human intents for collaboration. We detail in Section 2 the key design decisions we made to ensure our system is scalable and interpretable.

**2. We perform a comprehensive evaluation of MOSAIC.** We extensively test the limits of our system through 60 end-to-end trials where two robots collaborate with a human user to cook complex, long-horizon recipes (Section 3.2). We also test individual modules with 180 episodes of visuomotor picking (Section 3.4), 60 episodes of human motion forecasting (Section 3.5), and 46 online user evaluations of the task planner (Section 3.3)

**3. We analyze both MOSAIC's successes and failures to derive actionable insights for the field.** The modular nature of our architecture lends itself to error diagnosis — for each failure, one can clearly pinpoint which component failed and why it did. To this end, we distill key findings from our evaluations into limitations and exciting directions of future work based on our current system.

## 2 Approach

We present MOSAIC, **Mo**dular **S**ystem for **A**ssistive and **I**nteractive **C**ooking, a modular architecture that combines multiple large-scale pre-trained models to solve collaborative cooking tasks. Fig. 2 shows the three main components of MOSAIC: *1) Interactive Task Planner (2.1):* a module that interacts with real users via natural language to plan a diverse set of tasks and coordinate subtasks during the cooking process; *2) Human Motion Forecasting (2.2):* a module that leverages motion forecasting models to predict human motion such that robots can seamlessly collaborate with humans while maintaining a margin of safety to avoid human-robot collisions; *3) Visuomotor Skill (2.3):* a module that generalizes robot skills to a diverse set of kitchen objects and environments.

We make a set of simplifying assupmtions in our work, which we detail in Appendix B.

### 2.1 Interactive Task Planner

The goal of the task planner is to continuously interact with a human user using natural language, delegate subtasks to different robots or the user, and monitor progress. After agreeing upon a task with the user (e.g. "Prepare vegetable soup"), the task planner uses an LLM and an online recipe to represent the task as a direct acyclic graph (DAG) to model temporal dependencies between
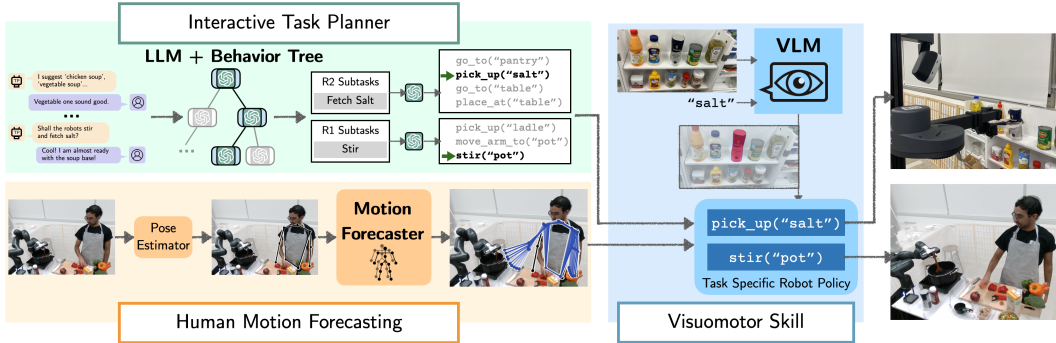
Figure 2: **MOSAIC System Overview.** The *Interactive Task Planner* module communicates with the user via natural language to decide on a recipe. It assigns subtasks to each robot accordingly. The *Human Motion Forecasting* module extracts and converts the human's 2D post to 3D coordinates, which it uses to predict future human motion. Separately, a VLM takes image and language as input and produces a 3D grasp pose around the object of interest. Combined, all three are taken by the execution policy of the *Visuomotor Skill* module to produce a final robot action.

different subtasks. From this representation, the task planner assigns and maintains a queue of subtasks for each robot based on the current agents' status and the user's requests. To execute a subtask (e.g. "fetch salt"), the task planner generates a code snippet that issues a series of API calls such as `go_to("pantry")`, `pick("pepper")`, etc. The initial set of subtasks and their mutual dependencies is generated once after the recipe is decided upon, the task planner can reassign subtasks and accept new ones based on the user's input.

While many recent approaches [1, 6, 7, 8, 9] directly use LLMs for task planning, we observe two main challenges. First, due to the large action space, the LLMs tend to violate constraints that the developer specifies even with chain-of-thought prompting [10]. Second, this approach requires all constraints to be specified in one monolithic prompt, which gives the developer little control over the LLMs' behavior and is challenging to debug and scale. To overcome both challenges, we propose an architecture that embeds LLMs within a behavior tree (BT) [11] (as shown in Figure 2). The nodes break the entire reasoning process down into easier reasoning process for LLMs to think about, thereby reducing the complexity and potential error rate of the LLMs and making it easy to scale to multiple behaviors. Finally, adding a new behavior is as simple as creating a prompt for that behavior and adding it as an option for other behaviors to invoke. No change to the code is necessary. Further details in Appendix C and exact prompts for all nodes are listed in Appendix I.

## 2.2 Human Motion Forecasting

Seamless and fluid coordination with humans while maintaining a margin of safety requires forecasting human motion. However, accurately forecasting human motion in dynamic environments such as kitchens is challenging, as humans can perform a wide range of motions, such as manipulating various objects in the kitchen or moving between stations. Even with large amounts of training data and a long context window, current state-of-the-art models struggle to accurately predict human motion at all times. We aim to build a forecasting model that generates predictions that sufficiently capture the impact of forecasted human motion during interactions with the robot.

**Training Pipeline.** We first pre-train our model on AMASS [12] a large dataset of human activity, encompassing over 300 subjects and 40 hours of motion capture data. However, AMASS only consists of general single-human movements (e.g. jumping, walking, dancing), which are not representative of a human's motion when working in close proximity with a robot. To this end, we utilize the Collaborative Manipulation Dataset (CoMaD) [13], a dataset consisting solely of human-human interaction episodes in a kitchen setting. For each episode in CoMaD, we identify these periods where both humans are in close proximity and construct a *transition dataset*. We sample data equally from the *transition dataset* and the entire CoMaD dataset to train the motion forecaster. This approach helps the forecaster maximize task efficiency by upsampling critical periods of the interaction when the human is likely to approach and interact with the robot.

**Inference Time: Real-time, Vision-based Forecasting and Planning.** Given an RGB-D image of the human, we use MediaPipe [14] to extract the 2D positions of their upper-body joints. These locations are then back-projected to 3D world coordinates using the image depth map and used to generate real-time motion forecasts for robot planning. However, the 3D coordinates obtained using the RGB-D camera are noisier than the high-fidelity motion capture data which the model is trained on. By injecting random Gaussian perturbations into the model's input at training time, we force the forecaster to learn to denoise potentially noisy input and generate smooth forecasts.

## 2.3 Visuomotor Skills

The visuomotor skills module takes the goal specification from the task planner, forecasts from the forecasting module, and observations from cameras and outputs a series of actions to complete the task. A common approach is end-to-end training on a suite of demonstrations [15, 16, 17, 18, 19, 20, 21], though good test-time performance generally requires optimal demonstrations to have good coverage over the observation space, which are challenging and time-consuming to collect. Instead, we decouple perception and action generation into a set of individual modules, where different skills have the flexibility to use different combinations as shown in Table 1. We describe these modules below, with more details in App. D.

| Skill | Freq. | Obj. Detect | Forecast | Action Exec. |
|---|---|---|---|---|
| pick(<obj>) | 31% | ✔ | | Learned |
| place(<loc>) | 20% | | | Engineered |
| stir(<obj>) | 11% | ✔ | ✔ | Engineered |
| handover() | 7% | | ✔ | Engineered |
| pour(<obj>) | 5.5% | ✔ | | Engineered |
| go_to(<loc>) | 25.5% | | | Planning |

Table 1: **Subcomponents Used per Skill.** pick() takes up the highest proportion of calls, and requires a policy that satisfies a tight set of constraints. Thus, we learn the policy via reinforcement learning (*RL*). go_to() assumes access to a map of the environment, and plans a path between the given start and goal.

**Object detection.** Given image and language input, we obtain a set of bounding boxes from OwlViT [22] and take the bounding-box coordinate with the highest CLIP similar score [23]. We use FastSAM [24] to obtain a more accurate segmentation of the object within the bounding box, and back-project the segmented pixels through the depth camera's point clouds and take the goal location to be its center of mass.

**Action execution.** For simple skills where there are no obstacles in the path between the robot and object-of-interest, we use an IK-based controller and a set of engineered primitives (e.g. a stirring motion), which take as input the object position and human motion forecasts to solve the task. For skills where such collisions may naturally occur during execution, we train an off-the-shelf reinforcement learning algorithm in a simulator that has access to an approximate robot dynamics model and a target grasp position. No action is executed in the simulator if the predicted action results in a collision, and instead gets a negative reward proportional to the distance from the goal. Further details on the motion primitives, simulator, and reward function are in Appendix D.

**Integration with human motion forecasts.** Some skills use the current end-effector position and user's motion forecasts (from Section E) to avoid collisions with the human (e.g. user drops food into the pot while the robot executes stir()) and anticipate future positions (e.g. handover()).

## 3 Experiments

### 3.1 Setup

In all experiments, the mobile manipulator is a 6-DoF Stretch Robot RE1 [25], and the tabletop manipulator is a 7-DoF Franka Emika Research 3 [26]. The kitchen has two overhead RGB-D cameras that can perceive the workspace and capture a human's motion. To allow users to interact with the task planner, we use Google's speech-to-text APIs [27] to transcribe user's verbal instructions and its text-to-speech APIs to vocalize the task planner's responses.

### 3.2 End-to-end Trials

The goal of end-to-end trials is to categorize how much failure each module has when the entire system is integrated and running together. To this end, we conduct a total of 60 end-to-end trials with

| Recipe | Typical Robot Skills Used | | Success | Subtasks Comp. |
|---|---|---|---|---|
| Toss Salad | R1: Go-to, Pick, Go-to, Place | R2: Pick, Stir | 8 / 10 | 92.5% |
| Tuna Sandwich | R1: Go-to, Pick, Go-to, Place, Handover, Go-to, Place | R2: Pick, Stir | 8 / 10 | 96.0% |
| Vegetable Soup | R1: Go-to, Pick, Go-to, Place | R2: Pick, Handover, Pick, Stir | 8 / 10 | 96.0% |
| Corn Soup | R1: Go-to, Pick, Go-to, Place | R2: Pick, Pour, Place, Pick, Stir | 6 / 10 | 90.0% |
| Caesar Salad | R1: Go-to, Pick, Go-to, Place | R2: Pick, Pour, Place, Pick, Stir | 5 / 10 | 86.7% |
| Chicken Soup | R1: Go-to, Pick, Go-to, Place, Handover, Go-to, Place | R2: Pick, Handover, Pick, Pour, Place, Pick, Stir | 6 / 10 | 91.4% |
| | | | 41 / 60 | 91.6% |

**Robot Subtasks:** Fetch sth, Handover sth, Stir sth, Put away sth, Pour sth

**Failure Cases:**
- (A) Pick Failed — 10
- (B) Place Failed — 3
- (C) Dropped Obj — 3
- (D) Interrupt Failed — 3
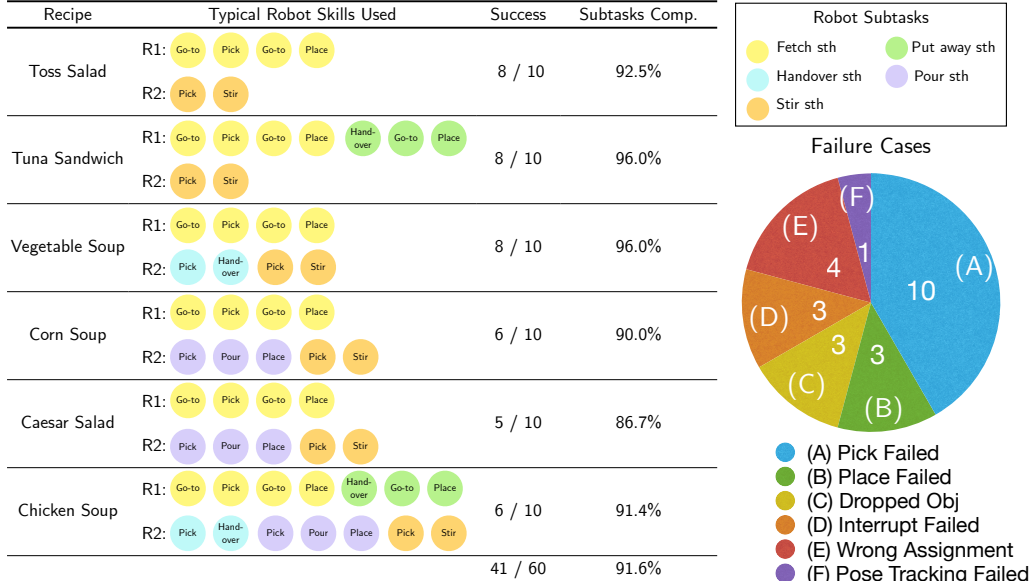- (E) Wrong Assignment — 4
- (F) Pose Tracking Failed — 1

Figure 3: **End-to-end results.** On-policy results for 6 recipes, where each recipe is tested through 10 trials. Each recipe contains various subtasks involving different robot skills. We report the number of trials that are completed without any errors and the individual subtask completion rate. We also categorize the failure cases. MOSAIC is able to complete 41/60 tasks with an average subtask completion rate of 91.6%.

two robots and a user collaboratively making 6 recipes. Each recipe involves a different combination of robot skills and different types of interaction with the user.[1] For example, users may provide vague instructions, interrupt a robot's subtask, and add new subtasks that are not in the recipe. Overall, MOSAIC completes 41/60 (68.3%) collaborative cooking trials of 6 different recipes with an average subtask completion rate of 91.6%.

Modularity enforces each sub-module to have a clear input/output contract, allowing one to *localize failures* and extract *transferable insights*. We use this to cluster failures into the 6 categories as shown in Figure 3. Specifically, errors originating from the task planner module usually come due to incorrect transcriptions of the user command from the text-to-speech sub-module. Errors in the perception module of visuomotor skills lead to an incorrect object identification or an insufficiently stable grasp. Likewise, tracking errors arise when the user moves outside the camera's view. At heart, *for complex tasks with multiple humans and robots, where failures are inevitable, modularity makes it easy to triage and treat failures*. In the following sections, we lift the insights from the end-to-end trials to each module and analyze how to limit errors therein.

### 3.3 Interactive Task Planner

Since the task planner directly interacts with the users, a frequent failure mode is constraint violation (e.g. acting without permission). Thus, it is crucial for it to exhibit predictable behaviors, especially when planning over long task horizons. To this end, we first quantify the frequency in which each constraint gets violated when the task planner interacts with real users in an online user[2]. We compare the proposed approach, which embeds LLMs within a behavior tree (*Tree*), against directly calling the LLM once with one prompt (*One-Prompt*). The monolithic LLM prompt has constraints that it must follow, explanations of what actions to choose in each situation, and in-context examples.

In the study, each user is randomly assigned to interact with either *Tree* or *One-Prompt* and is asked to analyze if the task planner changes subtasks without the user's permission (Act Without Permission), claimed robots can do subtasks beyond their capabilities (Lying), and did not respond to the user's instruction (Ignore User). We explicitly instruct participants to engage in non-nominal inter-

---

[1]The trials had authors acting as users and involved a total of 4 users between the ages of 20-30.

[2]The user study was approved by the Institutional Review Board at the University. See Appendix G for experimental setup, user study interface, and survey questions.
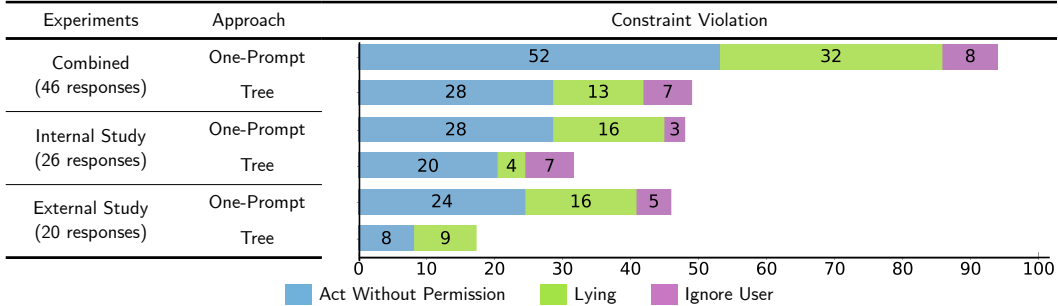
Figure 4: **Task Planner Constraint Violations During Real User Interactions.** We receive 46 responses in total (26 from internal and 20 from external study). Each user gets assigned either *Tree* or *One-Prompt* We present the total number of constraint violations per category. *Tree* makes 62.8% fewer constraint violations compared to *One-Prompt* for the combined responses, 36.2% fewer for internal, and 62.2% fewer for external.

actions with the task planner in hopes of eliciting a constraint violation. We received 26 responses from lab members who are not familiar with the task planner's capabilities and 20 responses from external users on Prolific [28], a crowd-sourcing website. Quantitatively, LLM with behavior tree violates 47.8% fewer constraints compared to the baseline of using one LLM prompt (Figure 4).

|  | Low Clutter | Med. Clutter | High Clutter | Total Success |
|---|---|---|---|---|
| OwlVit only [22] | 10/10 | 3/10 | 0/10 | 13/30 |
| OwlVit + CLIP (Ours) | 10/10 | 10/10 | 6/10 | 26/30 |

Table 2: **On-policy Evaluations of Different Vision Modules for `pick(<obj>)`** The architecture is tested on its ability to pick up the language-specified object when (i) a single object is in the pantry, (ii) 2-6 objects are in the pantry, and (iii) 7-15 objects in the pantry.

Table 7 in the Appendix provides examples and analysis of *One-Prompt*'s constraint violations. Finally, aggregated user feedback at the end of the survey indicates *Tree* generally provided a better user experience than *One-Prompt*. A user assigned with *Tree* stated "It worked as expected, quick and concise answers, compliant, didn't make any mistakes." Meanwhile, a user with *One-Prompt* commented "I could definitely see myself blowing my top with the level of disobedience."

Furthermore, localizing failure modes via modularity sufficiently scopes down the problem such that we can programmatically evaluate the approaches on unit tests. Specifically we test whether the task planner properly handles a request and chooses the right action as the interaction becomes more complex (e.g. user always disagrees with task planner and reassigns subtasks). *Tree* remains above 90.0% for its unit test pass rate, while *One-Prompt*'s performance drops from 100% (for easy cases) to 60.0% (for difficult cases), as the number of complex interactions increases. More details are in Table 9 in Appendix H. Overall, The results suggest that **compartmentalizing the action space by adding explicit structure to each LLM's reasoning problem significantly helps the task planner to respect constraints.**

### 3.4 Visuomotor Skills

End-to-end trials showed imprecision in the perception module to be a common source of failure for visuomotor skills. Thus, we now focus our analysis on quantifying how the vision component influences policy performance and qualifying common failure cases. We place the mobile manipulator in front of a pantry with increasing number of objects and test the success of `pick()`, the skill with the highest utilization frequency across all end-to-end runs.

Since localizing the error within the vision module, we found that directly using OwlViT [22] leads to rapid deterioration in accuracy when clutter increases, as the pre-trained model had difficulty identifying the correct bounding box from a large set of proposals. To rememdy this, we apply post-processing via Non-Maximum Suppression [29] and CLIP [23] (abbreviated as OwlViT + CLIP). Table 2 shows OwlViT + CLIP increases skill completion by 70% in medium-clutter regimes, and 60% in high-clutter regimes. Critically, **modularity helps to specifically identify the erring module, which allows improvements therein to have system-wide benefit.**

6

Additionally, figure 5 qualifies three failure cases of the vision module. First, the presence of too many objects (especially similarly-shaped ones such as various seasoning bottles) leads to a suboptimal set of bounding box proposals for CLIP to score. Second, when lighting and/or color blends the object contours into the background, only parts of the object may be included in the bounding box, resulting in a lower CLIP score. Lastly, imprecise prompts produce poor bounding box proposals.

### 3.5 Human Motion Forecasting

Finally, we analyze how using motion forecasts with noisy model inputs during training allow the module to be more robust during end-to-end trials through on-policy evaluations of a 7-DOF Franka robot arm collaborating with a real human user on two tasks.

The `stir()` skill involes the robot stirring a pot while the human periodically adds in vegetables. We measure the time it takes the robot to detect that the human arm is



Figure 5: **Vision backbone example failure cases.** We find that a cluttered background and poor lighting conditions to lead to a suboptimal set of bounding boxes for CLIP to score, while more specific prompts produce better bounding box proposals.

reaching into the pot (TIME TO REACT (MS)), the minimum distance maintained between robot arm and human hand (SAFETY MARGIN (CM)), and the number of times the human hand comes within a minimum threshold of the robot arm (COLLISIONS). In the `handover()` task, the user asks the robot to pick up and handover objects. We measure the average time to complete the handover (TIME TO GOAL (MS)) and the movement efficiency (PATH LENGTH (CM)), which measures the distance tracked by the robot's end-effector.

We compare against two baselines: (1) *Current* which assumes the current human pose will be its pose across the entire planning horizon, and (2) *Forecast (Base)* [30] that is not trained on noisy input. Each baseline and skill combination is evaluated 10 times for a total of 60 evaluations.

First, we find that using forecasts during the `stir()` skill significantly improves on all metrics, maintaining a 74% greater SAFETY MARGIN from the human on average and, more importantly, avoids any collisions. In contrast, we observe that the robot reacts very late when using the human's current pose, and results in collisions 20% of the time as shown in Table 3. In the `handover()` skill, the robot is 24% slower in completing the task following the current human wrist position compared to using the handover location predicted by our forecaster. Using the forecast, the robot moves directly toward the handover location, finishing the skill with 28% shorter trajectories.

Next, we ablate on training with noisy inputs by comparing our forecaster with a baseline approach (Base) [30] that does not train on noisy model inputs. In the `stir()` task, our forecaster has 23% quicker reaction time to human movements. Further, there is more variability in the performance of the Base forecaster (measured by the variance of each metric) that can be attributed to greater sensitivity to noisy inputs. Similarly, in the `handover()` skill, the Base forecaster's predictions are often erratic leading to jerky movements by the robot arm. Following them is no better than using the current human position for planning, as measured by task completion time and path length. Overall, we find **noise injections to the forecasting model make it more robust to perception errors, while using forecasts improves several key performance metrics of downstream skills.**

## 4 Discussion and Limitations

We decompose the overall problem of interactively cooking with a human user into a set of modules leveraging general-purpose pre-trained models. We localize errors in our system to individual modules and conduct targeted experiments. However, a number of open challenges still remain. First, improvements can be made to the task planner by grounding it with multi-modal input such as cameras and sensors. Further, expanding to new environments in a scalable and flexible manner may

| Task → | REACTIVE STIRRING | | | ROBOT TO HUMAN HANDOVERS | |
|---|---|---|---|---|---|
| Model ↓ | SAFETY MARGIN (cm) ↑ | TIME TO REACT (ms) ↓ | COLLISIONS ↓ | TIME TO GOAL (s) ↓ | PATH LENGTH (cm) ↓ |
| Current | 13.5 (±0.2) | 135.4 (±10.4) | 2/10 | 1.54 (±0.1) | 31.5 (±1.2) |
| Forecast (Base) [30] | 19.9 (±0.2) | 64.9 (±9.8) | 0/10 | 1.67 (±0.2) | 32.7 (±3.0) |
| Forecast (Ours) | 23.1 (±0.2) | 48.9 (±5.0) | 0/10 | 1.15 (±0.1) | 22.4 (±0.2) |

Table 3: **Task-Specific Performance Metrics.** We evaluate the robot's interactions with the human user on 2 collaborative manipulation tasks. Integrating forecasts into the robot's skills improves fluidity and increases safety margin across all metrics. We observe that relying on the current human pose during REACTIVE STIRRING is risky and results in collisions. ROBOT-HANDOVER tasks are more efficient using forecasting.



— Current Pose   — Forecasted Pose

Figure 6: **On-Policy Reactive Stirring.** (Left) **Current**: Using the human's current pose results in a delayed robot reaction and a collision once the human's hand enters the pot. (Right) **Forecast** Using the forecasted human position results in a smoother interaction and quicker reaction time, avoiding a collision.

require us to revisit previous assumptions and adopting new capabilities. Currently, the system's capabilities remain static after being deployed in an everyday user's household. An exciting area of future research is to continuously learn from real-time human feedback and interactions.

## 5   Related Work

**Home Robots.** Recent research efforts have attempted to provide robots with generalist capabilities to sufficiently adapt to home-like environments [1, 2, 3, 4, 31]. However, many of these works [3, 4] are limited to completing predefined tasks that don't require explicit task planning, e.g. picking a single item. Liu et al. [31] similarly tackle open-vocabulary navigation, but still sidesteps the challenge of a dynamic environment by assuming a static representation of the world after initialization. On the other hand, some works consider multi-arm/multi-robot planning for collaborative tasks [9, 32, 33, 34]. For example, Mandi et al. [9] significantly constrains human-robot collaboration by forcing the human to complete a specific task before the robot proceeds with its own task. In this paper, we aim to overcome these limitations by designing a system that enables multiple robots to fluidly collaborate alongside humans to perform a wide range of tasks.

**Specific Modules.** Our interactive task planner module is similar to work in single-robot settings with clearly defined language goals that generate a list of actions as the plan [8, 35, 36, 37, 38, 39] and synthesize code that calls robot action API [6, 40, 41, 42]; however, most of these works are non-interactive. In contrast to the most similar interactive work [8], we are solving a multi-agent task planning problem involving two robots and a user and communicating with the user to allocate tasks properly. Our visuomotor skills module is similar to the family of prior work [15, 19, 20, 43, 44, 18, 18, 44, 45, 46] that leverages VLMs for object identification within manipulation tasks. However, in contrast to prior work [47, 48, 42, 49, 45, 46], we train our action policy using reinforcement learning in simulation where affordances are provided by the VLM and constraints are inherent to the simulator. Collaborative manipulation tasks near humans necessitate human motion prediction, traditionally bypassed by assuming a static human [50, 51]. Advances in neural networks and the availability of extensive human motion datasets [12, 52, 13] have enabled the development of sophisticated RNN and GNN models to predict movements from past joint positions [53, 54, 55, 56, 57, 58, 30].

**Acknowledgments**

# References

[1] M. Ahn, A. Brohan, N. Brown, Y. Chebotar, O. Cortes, B. David, C. Finn, C. Fu, K. Gopalakrishnan, K. Hausman, et al. Do as i can, not as i say: Grounding language in robotic affordances. *arXiv preprint arXiv:2204.01691*, 2022.

[2] M. Bajracharya, J. Borders, R. Cheng, D. M. Helmick, L. Kaul, D. Kruse, J. Leichty, J. Ma, C. Matl, F. Michel, C. Papazov, J. Petersen, K. Shankar, and M. Tjersland. Demonstrating mobile manipulation in the wild: A metrics-driven approach. In K. E. Bekris, K. Hauser, S. L. Herbert, and J. Yu, editors, *Robotics: Science and Systems XIX, Daegu, Republic of Korea, July 10-14, 2023*, 2023. doi:10.15607/RSS.2023.XIX.055. URL https://doi.org/10.15607/RSS.2023.XIX.055.

[3] N. M. M. Shafiullah, A. Rai, H. Etukuru, Y. Liu, I. Misra, S. Chintala, and L. Pinto. On bringing robots home. *arXiv preprint arXiv:2311.16098*, 2023.

[4] S. Yenamandra, A. Ramachandran, K. Yadav, A. Wang, M. Khanna, T. Gervet, T.-Y. Yang, V. Jain, A. W. Clegg, J. Turner, et al. Homerobot: Open-vocabulary mobile manipulation. *arXiv preprint arXiv:2306.11565*, 2023.

[5] I. Idrees, T. Yun, N. Sharma, Y. Deng, N. Gopalan, G. Konidaris, and S. Tellex. Improved inference of human intent by combining plan recognition and language feedback, 2023.

[6] 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, 2023.

[7] B. Liu, Y. Jiang, X. Zhang, Q. Liu, S. Zhang, J. Biswas, and P. Stone. Llm+p: Empowering large language models with optimal planning proficiency, 2023.

[8] B. Li, P. Wu, P. Abbeel, and J. Malik. Interactive task planning with language models, 2023.

[9] Z. Mandi, S. Jain, and S. Song. Roco: Dialectic multi-robot collaboration with large language models, 2023.

[10] J. Wei, X. Wang, D. Schuurmans, M. Bosma, B. Ichter, F. Xia, E. Chi, Q. Le, and D. Zhou. Chain-of-thought prompting elicits reasoning in large language models, 2023.

[11] M. Colledanchise and P. Ögren. Behavior trees in robotics and AI: an introduction. *CoRR*, abs/1709.00084, 2017. URL http://arxiv.org/abs/1709.00084.

[12] N. Mahmood, N. Ghorbani, N. F. Troje, G. Pons-Moll, and M. J. Black. AMASS: Archive of motion capture as surface shapes. In *International Conference on Computer Vision*, pages 5442–5451, Oct. 2019.

[13] K. Kedia, A. Bhardwaj, P. Dan, and S. Choudhury. Interact: Transformer models for human intent prediction conditioned on robot actions. *ArXiv*, abs/2311.12943, 2023.

[14] V. Bazarevsky, I. Grishchenko, K. Raveendran, T. L. Zhu, F. Zhang, and M. Grundmann. Blazepose: On-device real-time body pose tracking. *ArXiv*, abs/2006.10204, 2020.

[15] 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.

[16] C. Chi, S. Feng, Y. Du, Z. Xu, E. Cousineau, B. Burchfiel, and S. Song. Diffusion policy: Visuomotor policy learning via action diffusion. *arXiv preprint arXiv:2303.04137*, 2023.

[17] E. Jang, A. Irpan, M. Khansari, D. Kappler, F. Ebert, C. Lynch, S. Levine, and C. Finn. Bc-z: Zero-shot task generalization with robotic imitation learning. In *Conference on Robot Learning*, pages 991–1002. PMLR, 2022.

[18] A. Stone, T. Xiao, Y. Lu, K. Gopalakrishnan, K.-H. Lee, Q. Vuong, P. Wohlhart, B. Zitkovich, F. Xia, C. Finn, et al. Open-world object manipulation using pre-trained vision-language models. *arXiv preprint arXiv:2303.00905*, 2023.

[19] A. Brohan, N. Brown, J. Carbajal, Y. Chebotar, X. Chen, K. Choromanski, T. Ding, D. Driess, A. Dubey, C. Finn, et al. Rt-2: Vision-language-action models transfer web knowledge to robotic control. *arXiv preprint arXiv:2307.15818*, 2023.

[20] M. Shridhar, L. Manuelli, and D. Fox. Cliport: What and where pathways for robotic manipulation. In *Conference on Robot Learning*, pages 894–906. PMLR, 2022.

[21] S. Nair, A. Rajeswaran, V. Kumar, C. Finn, and A. Gupta. R3m: A universal visual representation for robot manipulation. *arXiv preprint arXiv:2203.12601*, 2022.

[22] G. Heigold, M. Minderer, A. Gritsenko, A. Bewley, D. Keysers, M. Lučić, F. Yu, and T. Kipf. Video owl-vit: Temporally-consistent open-world localization in video. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 13802–13811, 2023.

[23] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.

[24] X. Zhao, W. Ding, Y. An, Y. Du, T. Yu, M. Li, M. Tang, and J. Wang. Fast segment anything. *arXiv preprint arXiv:2306.12156*, 2023.

[25] C. C. Kemp, A. Edsinger, H. M. Clever, and B. Matulevich. The design of stretch: A compact, lightweight mobile manipulator for indoor human environments, 2022.

[26] Franka research 3, 2022. URL https://franka.de/documents.

[27] URL https://cloud.google.com/speech-to-text/.

[28] Prolific, 2014. URL https://www.prolific.com.

[29] P. Viola and M. J. Jones. Robust real-time face detection. *International journal of computer vision*, 57:137–154, 2004.

[30] K. Kedia, P. Dan, A. Bhardwaj, and S. Choudhury. Manicast: Collaborative manipulation with cost-aware human forecasting. *ArXiv*, abs/2310.13258, 2023.

[31] P. Liu, Y. Orru, C. Paxton, N. M. M. Shafiullah, and L. Pinto. Ok-robot: What really matters in integrating open-knowledge models for robotics. *arXiv preprint arXiv:2401.12202*, 2024.

[32] M. Dogar, A. Spielberg, S. Baker, and D. Rus. Multi-robot grasp planning for sequential assembly operations. *Autonomous Robots*, 43:649–664, 2019.

[33] H. Ha, J. Xu, and S. Song. Learning a decentralized multi-arm motion planner. *arXiv preprint arXiv:2011.02608*, 2020.

[34] A. Tika and N. Bajcinca. Predictive control of cooperative robots sharing common workspace. *IEEE Transactions on Control Systems Technology*, 2023.

[35] M. Ahn, A. Brohan, N. Brown, Y. Chebotar, O. Cortes, B. David, C. Finn, C. Fu, K. Gopalakrishnan, K. Hausman, A. Herzog, D. Ho, J. Hsu, J. Ibarz, B. Ichter, A. Irpan, E. Jang, R. J. Ruano, K. Jeffrey, S. Jesmonth, N. J. Joshi, R. Julian, D. Kalashnikov, Y. Kuang, K.-H. Lee, S. Levine, Y. Lu, L. Luu, C. Parada, P. Pastor, J. Quiambao, K. Rao, J. Rettinghouse, D. Reyes, P. Sermanet, N. Sievers, C. Tan, A. Toshev, V. Vanhoucke, F. Xia, T. Xiao, P. Xu, S. Xu, M. Yan, and A. Zeng. Do as i can, not as i say: Grounding language in robotic affordances, 2022.

[36] W. Huang, F. Xia, T. Xiao, H. Chan, J. Liang, P. Florence, A. Zeng, J. Tompson, I. Mordatch, Y. Chebotar, P. Sermanet, N. Brown, T. Jackson, L. Luu, S. Levine, K. Hausman, and B. Ichter. Inner monologue: Embodied reasoning through planning with language models, 2022.

[37] W. Huang, P. Abbeel, D. Pathak, and I. Mordatch. Language models as zero-shot planners: Extracting actionable knowledge for embodied agents, 2022.

[38] K. Lin, C. Agia, T. Migimatsu, M. Pavone, and J. Bohg. Text2motion: from natural language instructions to feasible plans. *Autonomous Robots*, 47(8):1345–1365, Nov. 2023. ISSN 1573-7527. doi:10.1007/s10514-023-10131-7. URL http://dx.doi.org/10.1007/s10514-023-10131-7.

[39] S. Yao, D. Yu, J. Zhao, I. Shafran, T. Griffiths, Y. Cao, and K. Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. *Advances in Neural Information Processing Systems*, 36, 2024.

[40] 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, 2022.

[41] H. Wang, G. Gonzalez-Pumariega, Y. Sharma, and S. Choudhury. Demo2code: From summarizing demonstrations to synthesizing code via extended chain-of-thought, 2023.

[42] J. Wu, R. Antonova, A. Kan, M. Lepert, A. Zeng, S. Song, J. Bohg, S. Rusinkiewicz, and T. Funkhouser. Tidybot: personalized robot assistance with large language models. *Autonomous Robots*, 47(8):1087–1102, Nov. 2023. ISSN 1573-7527. doi:10.1007/s10514-023-10139-z. URL http://dx.doi.org/10.1007/s10514-023-10139-z.

[43] Y. Jiang, A. Gupta, Z. Zhang, G. Wang, Y. Dou, Y. Chen, L. Fei-Fei, A. Anandkumar, Y. Zhu, and L. Fan. Vima: General robot manipulation with multimodal prompts. *arXiv*, 2022.

[44] S. Karamcheti, S. Nair, A. S. Chen, T. Kollar, C. Finn, D. Sadigh, and P. Liang. Language-driven representation learning for robotics. *arXiv preprint arXiv:2302.12766*, 2023.

[45] F. Liu, K. Fang, P. Abbeel, and S. Levine. Moka: Open-vocabulary robotic manipulation through mark-based visual prompting. *arXiv preprint arXiv:2403.03174*, 2024.

[46] S. Nasiriany, F. Xia, W. Yu, T. Xiao, J. Liang, I. Dasgupta, A. Xie, D. Driess, A. Wahid, Z. Xu, et al. Pivot: Iterative visual prompting elicits actionable knowledge for vlms. *arXiv preprint arXiv:2402.07872*, 2024.

[47] 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.

[48] W. Huang, C. Wang, R. Zhang, Y. Li, J. Wu, and L. Fei-Fei. Voxposer: Composable 3d value maps for robotic manipulation with language models. *arXiv preprint arXiv:2307.05973*, 2023.

[49] P. Sundaresan, S. Belkhale, D. Sadigh, and J. Bohg. Kite: Keypoint-conditioned policies for semantic manipulation. *arXiv preprint arXiv:2306.16605*, 2023.

[50] W. Yang, B. Sundaralingam, C. Paxton, I. Akinola, Y.-W. Chao, M. Cakmak, and D. Fox. Model predictive control for fluid human-to-robot handovers. In *2022 International Conference on Robotics and Automation (ICRA)*, pages 6956–6962. IEEE, 2022.

[51] E. A. Sisbot and R. Alami. A human-aware manipulation planner. *IEEE Transactions on Robotics*, 28(5):1045–1057, 2012.

[52] C. Ionescu, D. Papava, V. Olaru, and C. Sminchisescu. Human3. 6m: Large scale datasets and predictive methods for 3d human sensing in natural environments. *IEEE transactions on pattern analysis and machine intelligence*, 36(7):1325–1339, 2013.

[53] H. Ling, G. Liu, L. Zhu, B. Huang, F. Lu, H. Wu, G. Tian, and Z. Ji. Motion planning combines human motion prediction for human-robot cooperation. In *2022 12th International Conference on CYBER Technology in Automation, Control, and Intelligent Systems (CYBER)*, pages 672–677. IEEE, 2022.

[54] V. Unhelkar, P. A. Lasota, Q. Tyroller, R.-D. Buhai, L. Marceau, B. Deml, and J. A. Shah. Human-aware robotic assistant for collaborative assembly: Integrating human motion prediction with planning in time. *IEEE Robotics and Automation Letters*, 3:2394–2401, 2018.

[55] J. Mainprice, R. Hayne, and D. Berenson. Predicting human reaching motion in collaborative tasks using inverse optimal control and iterative re-planning. *2015 IEEE International Conference on Robotics and Automation (ICRA)*, pages 885–892, 2015.

[56] V. Prasad, D. Koert, R. M. Stock-Homburg, J. Peters, and G. Chalvatzaki. Mild: Multimodal interactive latent dynamics for learning human-robot interaction. *2022 IEEE-RAS 21st International Conference on Humanoid Robots (Humanoids)*, pages 472–479, 2022.

[57] W. Mao, M. Liu, and M. Salzmann. History repeats itself: Human motion prediction via motion attention. In *European Conference on Computer Vision*, 2020.

[58] T. Sofianos, A. Sampieri, L. Franco, and F. Galasso. Space-time-separable graph convolutional network for pose forecasting. *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 11189–11198, 2021.

[59] J. X. Liu, Z. Yang, I. Idrees, S. Liang, B. Schornstein, S. Tellex, and A. Shah. Grounding complex natural language commands for temporal tasks in unseen environments, 2023.

[60] S. S. Raman, V. Cohen, D. Paulius, I. Idrees, E. Rosen, R. Mooney, and S. Tellex. Cape: Corrective actions from precondition errors using large language models. *arXiv preprint arXiv:2211.09935*, 2022.

[61] G. Brewka, T. Eiter, and M. Truszczyński. Answer set programming at a glance. *Commun. ACM*, 54(12):92–103, dec 2011. ISSN 0001-0782. doi:10.1145/2043174.2043195. URL https://doi.org/10.1145/2043174.2043195.

[62] Y. Jiang, S. Zhang, P. Khandelwal, and P. Stone. An empirical comparison of pddl-based and asp-based task planners. *CoRR*, abs/1804.08229, 2018. URL http://arxiv.org/abs/1804.08229.

[63] V. Lifschitz. Answer set programming and plan generation. *Artificial Intelligence*, 138(1):39–54, 2002. ISSN 0004-3702. doi:https://doi.org/10.1016/S0004-3702(02)00186-8. URL https://www.sciencedirect.com/science/article/pii/S0004370202001868. Knowledge Representation and Logic Programming.

[64] M. Fox and D. Long. PDDL2.1: an extension to PDDL for expressing temporal planning domains. *CoRR*, abs/1106.4561, 2011. URL http://arxiv.org/abs/1106.4561.

[65] A. Mavrogiannis, C. Mavrogiannis, and Y. Aloimonos. Cook2ltl: Translating cooking recipes to ltl formulae using large language models. *ArXiv*, abs/2310.00163, 2023.

[66] S. Y. Gadre, M. Wortsman, G. Ilharco, L. Schmidt, and S. Song. Clip on wheels: Zero-shot object navigation as object localization and exploration. *arXiv preprint arXiv:2203.10421*, 3 (4):7, 2022.

[67] A. Padalkar, A. Pooley, A. Jain, A. Bewley, A. Herzog, A. Irpan, A. Khazatsky, A. Rai, A. Singh, A. Brohan, et al. Open x-embodiment: Robotic learning datasets and rt-x models. *arXiv preprint arXiv:2310.08864*, 2023.

[68] S. Dasari, F. Ebert, S. Tian, S. Nair, B. Bucher, K. Schmeckpeper, S. Singh, S. Levine, and C. Finn. Robonet: Large-scale multi-robot learning. *arXiv preprint arXiv:1910.11215*, 2019.

[69] M. Minderer, A. Gritsenko, and N. Houlsby. Scaling open-vocabulary object detection. *arXiv preprint arXiv:2306.09683*, 2023.

[70] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.

[71] A. Raffin, A. Hill, M. Ernestus, A. Gleave, A. Kanervisto, and N. Dormann. Stable baselines3, 2019.

We first enumerate more related works in Section A, then outline the configuration of our system, as detailed in Section B. Subsequently, we delve into further discussions on each component of our system, covered across Sections C, D, and E. Finally, we present an in-depth analysis of the user study, incorporating details on the experimental setup and supplementary findings, all of which are elaborated in Section G.

## Appendix A  Extended Related Works

**Modular Architectures.** Modularization has been extensively employed to partition complex and long-horizon robotics tasks into more easily addressed sub-components. For example, in the space of task planning, different Large Language Model (LLM) modules are used to improve command interpretation in unseen environments [59] and to produce corrective replanning prompts [60]. OK-Robot [31] focuses on the problem of object retrieval and navigation using separate submodules for mapping, object-detection, path-planning, and grasping. Our approach tackles the combined domains of task planning, visuomotor skill learning, and human motion forecasting. We leverage modularity at multiple levels, e.g. we outperform a single VLM object detection by combining Owl-Vit to detect bounding boxes, and CLIP to select the correct box.

**Task Planning.** A task planner takes as input a high-level task, e.g. cooking a recipe, and generates a plan, e.g. a sequence of sub-tasks, to achieve that goal. Traditional approaches frame this as a search problem and invoke a symbolic planner to solve it [61, 62, 63, 64]. However, using these methods for everyday tasks is challenging because they require pre-defining the search space and lack a natural-language interface to interactively communicate the task. Recent work leverages LLMs for task planning to overcome both of these limitations. In single-robot settings, given a clearly defined language goal, recent work can be categorized as generating a list of actions as the plan [35, 36, 37, 38], synthesizing code that calls robot action API [6, 40, 41, 42], or translating to a problem solvable by a classical planner [7, 65]. However, none of these systems interact with humans and coordinate tasks for both humans and robots.

Idrees et al. [5] focuses on estimating a user's intent and using question-answering interaction with the user to update that estimate. Then, the robot can use that estimate to suggest how the user can make task progress. However, its interaction is limited to yes-or-no questions, and it does not focus on task planning for multiple agents based on the robot's capabilities and the user's unstructured natural language feedback.

Li et al. [8] has the closest task planning framework to our approach, where the LLM takes a specific natural language goal to generate a step-by-step plan before synthesizing robot code for each step. However, because we are solving a multi-agent task planning problem involving two robots and a user, our task planner cannot simply output a list of steps. It must continuously communicate with the user to properly allocate subtasks to suitable agents.

**Visuomotor Skills.** Several recent works study the application of pre-trained vision-language models (VLMs) to robotics [18, 66, 15, 19, 47, 20, 43, 44]. One family of recent work [15, 19, 20, 43, 44, 18] integrate pre-trained VLMs in an end-to-end fashion, e.g. segmenting out regions of interest to assist in action prediction [18, 44]. A second flavor of approach [47, 48, 1, 37] leverages VLMs to recognize affordances and constraints in the environment and provide corresponding execution instructions through language [47] or code [48]. Our model is similar in this aspect, where we distinguish the training objectives of environment perception and action execution. This effectively liberates us from needing a large dataset of humans or robots demonstrations to provide good coverage [18, 15, 19, 20, 21] and from having to worry about embodiment mismatch between large-scale robot learning datasets [67, 68].

**Human Motion Forecasting.** Collaborative manipulation tasks in close proximity to humans require predicting human motion. This is a challenging problem since human motion is complex and highly variable. A common approach is to sidestep the problem of motion forecasting [50, 51] by considering the human to be static. Instead, recent research is moving towards the use of neural

networks and supervised learning to predict future human motion based on a short history of past joint positions [53, 54, 55, 56]. The release of large open-sourced datasets of human motion [12, 52] has made it possible to train large RNN and GNN-based neural network models for human-pose forecasting [57, 58]. Consequently, these datasets have been integrated into robot motion planning, focusing on collaborative manipulation tasks [30, 13]. Closest to this work, ManiCast [30] proposed a framework to learn cost-aware human forecasts. However, this approach relies on a bulky motion camera setup, requiring the user to wear a motion capture suit with markers. In this work, we run our integrated human motion forecasting and planning system in real-time using a single RGB-D camera to track human pose.

## Appendix B  System Setup

**Kitchen Scene and Robot Placement.** The kitchen scene consists of a main kitchen table at the center where all cooking activities are performed. A pantry is placed near the table, which contains a large range of condiments and kitchen staples. There is also a secondary table on the side of the center table meant for serving up the final dishes. Our robot system includes two robots (R1 and R2).

- R1 (*Franka Emika Research 3* [26]) is a tabletop 7-of manipulator stationed at one end of the kitchen tables at the center of the scene.
- R2 (*Hello Robot Stretch RE1* [25] is a mobile manipulator that can navigate around the kitchen area, capable of fetching and putting away condiments and kitchenware as required by the user.

**Camera Placement.** For the tabletop manipulator (R1), the perception stack includes two Intel Realsense D435i RGB-D cameras placed above the center kitchen table. Both cameras are placed at opposite ends of the table and at an angle such that they capture the entirety of the tabletop as well as the human user. Integrating both camera perspectives enhances the visibility of objects and human poses within a cluttered kitchen setting, effectively mitigating occlusion issues. The mobile manipulator (R2) uses an onboard Intel Realsense D435i RGB-D head-camera for perceiving objects.

**Computational Details.** In addition to the onboard computing capabilities of the robots, our setup includes five personal computers (PCs) dedicated to running various system modules. These PCs are connected to the same network, utilizing the Robot Operating System (ROS) for communication. For tasks that demand real-time neural network inference, we employ onboard GPUs, (NVIDIA GeForce RTX 3060). Detailed information about each PC's role and configuration is provided:

- **C1:** Connected with a Bluetooth microphone and speaker, this PC runs the *Speech-Text* system for communicating with the user and the *Interactive Task Planner* that utilizes GPT-4 API calls.
- **C2:** Used for running neural network models related to the perception (object detection) and control (RL agent) of R2. This PC also communicates with C1 to allocate subtasks to R2.
- **C3:** This PC forms the perception stack for R1, including running neural network models for object detection and pose estimation.
- **C4:** This PC runs the human forecasting model using the pose estimates computed by C3. This PC also computes motion plans for R1 based on the predicted object pose and human forecasts. Further, it communicates with C1 to allocate subtasks to the R2.
- **C5:** This PC is installed with a real-time kernel to send joint commands to R2 at 1 kHz frequency as recommended by the robot manufacturers.

**System Assumptions.** In accordance with the setup above, we make a set simplifying assumptions in our work:

1. *Access to a set of seed recipes:* A recipe contains a set of subtasks with temporal dependencies. We seed the system with an initial set of recipes, but the user has the freedom to make modifications on the fly (e.g. adding an ingredient).

2. *Access to a map:* We assume that our system has mapped the kitchen ahead of time, so it is aware of where ingredients and tools are stored and how to navigate to different locations.

3. *Full observability:* We assume that objects are not occluded for detection and grasping, though they can be next to each other. We also assume that the upper torso of the human is visible to the cameras for tracking and prediction.

4. *Skills API:* We assume access to a library of robot skills that can be invoked with specific input parameters (e.g. `pick_up("salt")`, `stir()`).

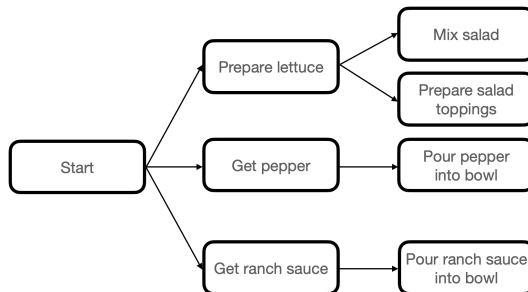## Appendix C  Interactive Task Planner Details



Figure 7: **Recipe DAG Example.** This DAG represents the subtasks and dependencies involved in making a Caesar Salad. At the beginning of making the recipe, the available subtasks include 'Prepare lettuce', 'Get pepper', and 'Get ranch sauce'. If one of the subtasks is marked as done, the following subtasks become available (e.g. completion of 'Get pepper' causes 'Pour pepper into bowl' to become available).

The interactive task planner consists of three main components: a representation of a task and the dependencies of its subtasks, a mechanism to decide on a recipe and assign subtasks to others, and a medium to communicate to robots which skills to use for a given subtask. We implement these using a direct acyclic graph (DAG), a behavior tree, and LLM-generated code (communicated over ROS action services).

**Task DAG.** The task planner represents a task (e.g. "Prepare vegetable soup") as a DAG, whose nodes represent subtasks of that task and whose edges represent dependencies between the subtasks. However, a DAG alone is insufficient for generating a task plan based on the user's instruction because it does not specify exactly which robot should complete which subtask.

For a recipe, the DAG is generated by an LLM prompt ahead of time. Concretely, the LLM takes as input the ingredient list and step-by-step instructions, scraped from the recipe website.[3] It outputs a marked-down nested list of the recipe, which is easier for the LLM to reason about and generate in practice. A list item represents a subtask, and the nested structure represents the dependencies, so we can programmatically convert the LLM output into a DAG. The complete LLM prompt is in the supplementary material.

The task planner also maintains a done state for each node/subtask in the DAG. To determine the available subtasks, we start from a root node whose done state is set, with outgoing edges to the first subtasks. Then, we follow each outgoing edge until reaching a node whose done state is unset and add it to a set. If no node is found through this process, the recipe has been finished. See Figure 7 for an example of a DAG for Caesar salad.

---

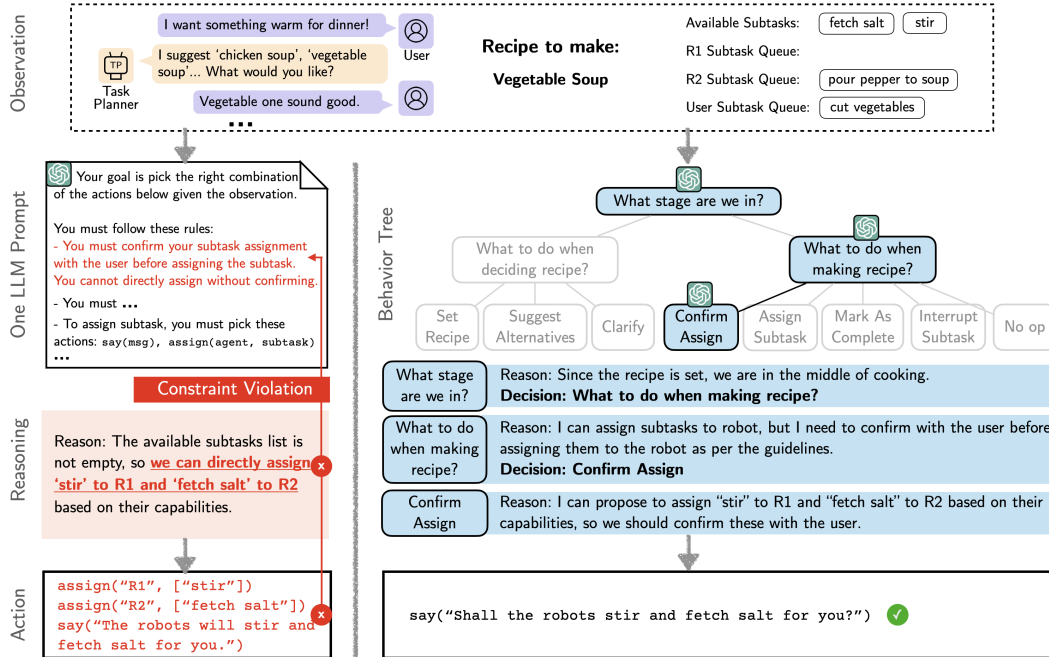[3]We used an existing open-source recipe scraper.

Figure 8: **Tree-structured task planner vs single-prompt LLM.** We compare our approach against using one LLM prompt, which tends to violate constraints. Given the observation, the LLM with one monolithic prompt directly assigns subtasks to robots, which violates the constraint that it must confirm with the human before assigning tasks. Meanwhile, because our approach compartmentalizes the action space and reasoning process in a behavior tree, it is able to follow a correct reasoning path and correctly confirm its subtask proposal with the user.

A DAG allows us to represent dependencies, such as, *sequential:* 'do A before B', and *AND* dependencies: 'do A and B before C'; it currently does not allow *OR* dependencies (do A or B before C). However, we could still create 32 unique recipes with this limitation.

**Behavior Tree.** We use a behavior tree to decide on a recipe and then assign subtasks to others by designing the tree around the behaviors we expect. Fig. 8 visualizes the behavior tree and shows an example comparing MOSAIC's task planner to the baseline that relies on a monolithic LLM prompt. Each behavior is encapsulated in a node, which represents a call to an LLM with a specific prompt and a pre-defined set of decisions to choose from. It takes as input an observation from the world and outputs arguments for high-level actions. For example, Fig. 9 shows snippets of prompts for different behaviors. The instructions describe the goal, the action space, which part of the observation to focus on, constraints to adhere to, and in-context examples. `Assign Subtask` is a leaf node that directly assigns subtasks to the robot and speaks to the user. On the other hand, the `What to do when making recipe?` behavior is



Figure 9: **Behavior node snippets.** Two prompt snippets of behavior nodes in our behavior tree. The top box shows a node that predicts a set of actions $a_t^{\text{high}}$ to execute. The bottom box shows a node that predicts which child node $n'$ to go to.

a higher-level node that calls other behaviors, e.g. `Confirm Assign`, `Assign Subtask`, etc.

More specifically, each sample from the observation space consists of:

1. recipe name

2. available subtasks

3. each robot's subtask queue, current subtask, and current status (`Idle`, `Running`, or `Interrupted`)

4. user's subtask queue

5. completed subtask queue

6. user's current input

7. chat history

The recipe name can be empty if the recipe has not been decided yet. The available subtasks are populated by the DAG. The robot and user subtasks are all populated by the behavior tree's high-level actions; the robot additionally has a current subtask and status field updated over a ROS action server as the robots complete their subtasks. When subtasks are completed, the completed subtask queue is updated. Finally, the user's currently spoken input is stored and later appended to the chat history along with the task planner's messages.

The high-level actions include

- `say(msg)`
- `set_recipe(name)`
- `assign(agent, subtasks)`
- `mark_complete(subtasks)`
- `interrupt(agent)`
- `no_op()`

`say(...)` allows the task planner to communicate to the user with a message. `assign(...)` will assign a list of subtasks to an agent (robot or human). `mark_complete(...)` will set a list of subtasks as completed. `interrupt(...)` will stop a robot from doing its current subtask. `no_op()` does nothing.

The tree consists of various nodes that each query an LLM that either outputs (1) a decision for the next node to run or (2) arguments for the high-level actions to take. Each node is associated with a prompt that is used when querying the LLM. If a node's query response is malformed (e.g. bad JSON) or invalid (e.g. bad decision or arguments), the node is rerun. Each node only requires the observation as input, so we can run each node simultaneously to parallelize the LLM queries and draw a path from the root to a leaf based on the decisions made.

The tree runs a cycle to take high-level actions whenever the observation differs from the past observation. This gives the user time to respond to the task planner's questions. To receive user input and respond to the user, we use speech-to-text and text-to-speech modules, respectively. The tree runs indefinitely until the script is terminated.

**Code Generation.** Whenever the task planner assigns a subtask to a robot, it must be converted into a sequence of low-level skills the robot is capable of. We do this by using an LLM to generate code that the robot runs.

When the task planner assigns a subtask to a robot, it is first added to the robot's subtask queue. A thread dedicated to the robot checks to see if there are any subtasks in the queue, pops it to add to the current subtask, and sets the status to Running. A separate prompt for code generation is used to query an LLM to generate code for the provided subtask. An example of generated code includes

```python
from robot_utils import <robot_api>
from env_utils import <env_constants>
pick_up_item(LADLE)
place_item_at(POT)
stir()
```

where `<robot_api>` includes all low-level robot skills like `pick_up_item(...)` and where `<env_constants>` includes enums for objects in the environment. Each line of code executing a robot skill sends a ROS action to the robot to execute said skill. When the robot finishes executing its current skill, it communicates that it has finished to the task planner, which can, in turn, send another skill. This continues until the entire subtask is finished, in which case the robot's current subtask is cleared, and the robot's status is set to Idle. If the robot is interrupted, its current subtask is also cleared, but its status is set to Interrupted.

## Appendix D    Visuomotor Skills Details

**Skill Library.** The task planner has access to a number of robot skills represented as function calls that are parameterized by object positions and target locations. For each skill, the positions of the objects are estimated using an open-vocabulary object detection model, OWL-ViT (more details in the next section), given text prompts provided by the task planner. For navigation, we store mapped locations to real-world coordinates, assuming the kitchen scene does not change its configuration between runs.

We enumerate below the set of low-level skills performed by the two robots in this paper:

1. `pick(<obj>)`: Both robots share the same object detection module to complete the `pick(<obj>)` task to get bounding boxes and a 3D grasp-pose around the object of interest. R1 (Franka arm) moves directly to the grasp pose using an inverse kinematics-based joint impedance controller. R2 (Stretch robot) is tasked with picking up objects from a cluttered pantry. To avoid hitting the pantry and surrounding objects, the robot uses a reinforcement learning policy trained in simulation to execute actions.

2. `go_to(<loc>)`: This skill uses a map of the kitchen acquired beforehand and the internal localization mechanism of Stretch RE1 to navigate to designated locations around the kitchen.

3. `place(<loc>)`: The `place(<loc>)` skill is parameterized by the target locations and completed with pre-coded motion primitives.

4. `stir(<obj>)` We define this motion primitive for R2 (Franka arm) holding a tool (such as a ladle) in its arm parameterized by the target utensil where the action takes place (`<obj>`), for example, a pot. Further, this skill is responsive to the human's movements in the robot's stirring radius. If the human's motion forecasts reach into the robot's workspace, the robot stops stirring and makes space for the human to move in.

5. `pour(<obj>)` Similar to the `stir()` function, this skill enables R2 to pour an already gripped object such as a salt can into a target receptacle (`<obj>`), such as a bowl. This process involves the utilization of motion primitives based on the estimated locations of the objects involved. Specifically, in the scenario of pouring salt into a bowl, R2 executes a sequence of actions: it first positions the salt can over the bowl at a calculated tilt angle and then shakes the can to dispense the salt. Following the completion of the pouring action, R2 returns the salt can to its original location on the table.

6. `handover()` R2 (Franka arm) completes handovers quickly and efficiently by directly moving its end-effector towards the forecasted human wrist position. Once the robot's end effector is within a threshold of the human's wrist position, it stops and releases the object into the robot's hand. Finally, the robot arms reset back to its original position.

**Object Localization.** The object localization pipeline first takes as input RGB image and text prompt of the object of interest, which is passed through an OWLViT [69] object detection model that produces $k$ bounding box proposals denoting possible locations of the object. These $k$ bounding boxes are then filtered using non-maximum suppression to remove overlapping boxes. Due to the camera angle and other noise in the environment, we find that the top OWLViT bounding box does not reliably agree with the desired object. Thus, these proposals are refined by feeding each of the

images of the cropped bounding boxes and the text prompt to a pre-trained CLIP [23] model to create a CLIP score that measures how aligned each cropped image is with the text prompt[4].

Next, the image, the bounding box with the highest CLIP score, and the text prompt is fed to a pre-trained FastSAM [24] model to segment the object located in the bounding box. The point cloud given by the depth camera is used to project all the points inside the segmentation mask into 3D space. All the 3D points of the object are averaged to obtain a final, single 3D point. This 3D point is then fed to the execution module to produce actions for how the robot should move to the object.

**RL Simulator and Reward Function.**

The RL agent needs to take the goal prediction and execute a series of actions to reach that goal without collisions. For `pick()` specifically, consider a pantry that is stocked with items. A desirable trajectory would avoid hitting the pantry boards, hitting neighboring objects, and pushing the object as the gripper approaches. To guide the agent, we create a simulator that, for a given goal point, builds a 3-dimensional set of walls to the sides, back, and bottom of the goal. Invalid actions are those that collide with a wall or violate robot joint states. An episode starts by sampling a start and goal position within some distance reachable by our robot.
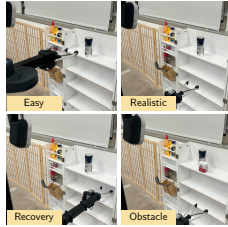
The observation space is the $L_1$ norm between the goal and current positions. We then train a Proximal Policy Optimization [70] agent using the implementation from Raffin et al. [71] with the same action space as the teleoperation commands in the demonstration data using the following cost function

$$\exp(-\|O_c - O_g\|_2) - 1 \qquad (1)$$

where $O_c$ and $O_g$ represent the current and desired end-effector coordinates respectively, and $\| \cdot \|_2$ is the Euclidean distance. The main failure case for the agent is violating joint constraints while trying to avoid the walls because the observation space does not include joint states.



|  | Easy | Realistic | Recovery | Obstacle |
|---|---|---|---|---|
| RL | 10/10 | 10/10 | 10/10 | 4/10 |
| IK | 10/10 | 3/10 | 0/10 | 2/10 |
| BC | 6/10 | 4/10 | 0/10 | 3/10 |

Table 4: **On-policy Evaluations of Policy Module.** We evaluate under four different starting configurations: a) *Easy*, when the gripper is close to the object; b) *Realistic*, when the gripper is retracted; c) *Recovery*, when the gripper is in an extended position away from the object; and d) *Obstacle*, when the object is partially occluded. We see the RL agent trained in simulation successfully reaches the goal without hitting the pantry, despite being reset to states that oblige recovery motions. However, success rate deteriorates as object placements violate the initial assumptions made about the simulator used to train the agent.

**Comparing Training RL in Simulation with BC and IK.** We evaluate the action-execution module of the `pick()` skill from 4 different types of starting configurations, visualized in Table 4. The RL agent completes the skill and avoids collision, achieving the highest success rate of 85% across the different configurations.

Using an IK-based controller successfully reaches the goal in the *Easy* configuration where the gripper is directly in front of the object. In other harder settings, such an approach often collides with the pantry, leading to an overall success rate of $32.5\%$. We also train a behavior cloning (BC) agent using 50 demonstrations of the `pick()` across the *Easy* and *Realistic* configuration. BC achieves $50\%$ accuracy when tested within these two settings, but its performance decreases to $15\%$ in out-of-distribution configurations (*Recovery* and *Obstacle*). The RL policy has a perfect success rate across all configurations except for *Obstacle*, which it completes 40% of the time. This is because we assumed an absence of occlusion while designing our reward function for the RL policy, and thus the *Obstacle* configuration demonstrates a limitation to learning via hand-designed reward functions. We posit that with sufficient demonstrations, BC is capable of learning a more expressive policy in such situations.

---

[4]If $k$ is set too low, the set may not contain a bounding box around the object of interest to be used by CLIP. If $k$ is set too high, the set of bounding boxes may be too noisy, resulting in lower accuracy. We set $k = 10$ for all experiments.

**Behavioral Cloning Baseline** Our BC policy consists of two feed-forward layers with 256 neurons and is trained on 50 demonstration trajectories with variation in the robot arm's starting height and location of the object. At each timestep, the model takes as input the difference between the current end-effector position and the final position (the same as the RL agent), and outputs a 10-dimensional vector of logits, where each dimension corresponds to moving one of the robot's 10 joints. The model is trained using a weighted cross-entropy loss function to account for class imbalances. On-policy, a final action is obtained by categorically sampling from the output vector.

## Appendix E   Human Motion Forecasting

**Model Architecture.**   We use a Space-Time Separable Graph Convolutional Network (STS-GCN) [58] model architecture for our human-motion forecaster, which encodes the human's joint positions at different timesteps as nodes in a graph. Instead of simply constructing a fully connected graph between all nodes, the model constructs a sparse network without redundant edges across temporal and spatial dimensions. Edges are connected only between the same human joint through consecutive timesteps and between all joints at the same timestep.

**Experimental Setup.** In order to employ our human motion forecasting model for real-time inference, we make use of an RGB-D camera (Intel RealSense D435) pointed at the human's torso. The human pose is represented by the 3D positions of 7 upper body joints (shoulders, elbows, wrists, and neck). We track the 2D human joint locations using MediaPipe [14] on input RGB images and back-project them to 3D world coordinates using the depth map. As discussed in the approach, our method is forced to handle noisy inputs from depth map projections which are out-of-distribution for motion forecasting models trained on high-fidelity motion capture data. We first compare forecasting performance on CoMaD [13] to select a model suitable for predicting human motion in our dynamic kitchen setting. CoMaD is collected via motion capture suits and contains 270 episodes of human-human interactions across 3 different kitchen tasks with an average length of 30 seconds per episode (4+ hours of total data). Then, we conduct experiments injecting various levels of random Gaussian noise into motion capture data at train time to overcome the train-test distribution mismatch and report results on a dataset of human motion tracked by our single-camera setup.

**Forecasting Metrics.** We quantify errors made by the forecaster by measuring both the Average Displacement Error (ADE) on all predicted timesteps and the Final Displacement Error (FDE) of the predicted pose 1-second prediction into the future given 0.4 seconds of pose history. We report metrics on All Joints as well as Wrists specifically, as they are the most relevant joints in the manipulation tasks we roll out. We additionally report forecasting metrics on the CoMaD *transition dataset* of human motion during short transition windows in which humans come in close contact with one another, denoted by prefix 'T-' (e.g. T-All Joints ADE, T-Wrists ADE). Note that humans are always in very close proximity during the TABLE SETTING task.

**CoMaD Forecasting Results.** Our two baselines are (1) BASE, trained only on AMASS data, and (2) SCRATCH, trained only on CoMaD data. We report results for two more models: (3) FINE-TUNED, pre-trained on AMASS data and fine-tuned on CoMaD data, and (4) FINETUNED-T, pre-trained on AMASS data and fine-tuned on CoMaD with upsampling from its *transition dataset*. Each model is tested on a held-out CoMaD test set of episodes. FINETUNED-T significantly outperforms all other models across every metric for the REACTIVE STIRRING and HANDOVER tasks. On the TABLE SETTING task, FINETUNED only marginally produced lower errors compared to FINETUNED-T, both of which beat out the baselines.

We find that upsampling CoMaD transition data where humans are in close contact enables more accurate motion forecasts on kitchen activities. BASE struggles to generate accurate predictions in highly dynamic manipulation tasks, as it was only trained on AMASS [12] data of a single-human and lacks interaction data. SCRATCH is challenged with learning general human motion dynamics from CoMaD, a much smaller dataset compared to AMASS, which is reflected by its higher errors. Ultimately, we find that pre-training forecasting models on large-scale human activity data and fine-

| Metrics (mm) ↓ | | BASE | SCRATCH | FINETUNED | FINETUNED-T |
|---|---|---|---|---|---|
| **REACTIVESTIR** | All Joints ADE | 60.3 (± 0.6) | 40.0 (± 0.3) | 32.1 (± 0.2) | 29.9 (± 0.2) |
| | All Joints FDE | 91.5 (± 0.9) | 60.3 (± 0.5) | 54.0 (± 0.5) | 51.7 (± 0.4) |
| | Wrists ADE | 83.7 (± 0.6) | 58.0 (± 0.4) | 47.9 (± 0.3) | 44.9 (± 0.3) |
| | Wrists FDE | 128.0 (± 1.0) | 87.2 (± 0.7) | 80.7 (± 0.6) | 76.6 (± 0.6) |
| | T-All Joints ADE | 58.0 (± 0.4) | 38.7 (± 0.2) | 31.1 (± 0.1) | 28.8 (± 0.1) |
| | T-All Joints FDE | 87.7 (± 0.6) | 58.0 (± 0.3) | 52.0 (± 0.3) | 49.6 (± 0.3) |
| | T-Wrists ADE | 81.8 (± 0.4) | 56.8 (± 0.2) | 46.8 (± 0.2) | 43.8 (± 0.2) |
| | T-Wrists FDE | 124.6 (± 0.7) | 84.9 (± 0.5) | 78.7 (± 0.4) | 74.4 (± 0.4) |
| **HANDOVER** | All Joints ADE | 56.3 (± 0.3) | 40.4 (± 0.2) | 32.9 (± 0.1) | 31.4 (± 0.1) |
| | All Joints FDE | 88.0 (± 0.5) | 62.8 (± 0.4) | 56.2 (± 0.3) | 55.0 (± 0.3) |
| | Wrists ADE | 88.5 (± 0.4) | 64.2 (± 0.3) | 51.8 (± 0.3) | 50.0 (± 0.2) |
| | Wrists FDE | 139.4 (± 0.8) | 100.3 (± 0.6) | 89.2 (± 0.6) | 87.4 (± 0.6) |
| | T-All Joints ADE | 54.0 (± 0.2) | 38.9 (± 0.1) | 31.7 (± 0.1) | 30.2 (± 0.1) |
| | T-All Joints FDE | 83.8 (± 0.4) | 59.6 (± 0.3) | 53.5 (± 0.3) | 52.4 (± 0.3) |
| | T-Wrists ADE | 85.2 (± 0.3) | 61.9 (± 0.3) | 50.1 (± 0.2) | 48.3 (± 0.2) |
| | T-Wrists FDE | 133.0 (± 0.6) | 95.4 (± 0.5) | 85.2 (± 0.4) | 83.4 (± 0.4) |
| **TABLESET** | All Joints ADE | 107.0 (± 1.1) | 72.0 (± 0.5) | 59.0 (± 0.4) | 59.1 (± 0.4) |
| | All Joints FDE | 181.0 (± 1.9) | 118.1 (± 0.9) | 108.0 (± 0.8) | 108.8 (± 0.8) |
| | Wrists ADE | 127.1 (± 1.0) | 93.4 (± 0.6) | 80.4 (± 0.5) | 81.7 (± 0.5) |
| | Wrists FDE | 224.7 (± 2.0) | 152.6 (± 1.1) | 143.1 (± 1.0) | 145.8 (± 1.0) |

Table 5: **CoMaD Forecasting Metrics.** We report Average Displacement Error (ADE) and Final Displacement Error (FDE) for Handover, Reactive Stirring, and Table Setting tasks on different forecasting models: Base, Scratch, FineTuned, and FineTuned-T. Metrics prefixed with 'T-' indicate measurements from the *transition dataset*, data where humans are in close-contact. Finetuned-T produces the lowest errors on Reactive Stirring and Handover, with very marginally higher errors on Table Setting.

| Metrics (mm) ↓ | | NOISE $_0$ | NOISE $_{0.001}$ | NOISE $_{0.01}$ | NOISE $_{0.1}$ |
|---|---|---|---|---|---|
| **REACTSTIR** | All Joints ADE | 75.1 (± 1.2) | 70.8 (± 1.2) | 64.8 (± 0.9) | 136.2 (± 0.9) |
| | All Joints FDE | 107.3 (± 1.8) | 103.5 (± 1.7) | 94.0 (± 1.3) | 155.4 (± 1.2) |
| | Wrists ADE | 97.6 (± 1.8) | 90.4 (± 1.8) | 81.8 (± 1.5) | 116.0 (± 1.3) |
| | Wrists FDE | 128.1 (± 2.5) | 124.5 (± 2.5) | 120.7 (± 2.1) | 140.3 (± 2.1) |
| **HANDOVER** | All Joints ADE | 66.1 (± 1.0) | 59.9 (± 1.0) | 55.2 (± 0.8) | 151.1 (± 0.5) |
| | All Joints FDE | 95.9 (± 1.4) | 90.6 (± 1.4) | 83.2 (± 1.2) | 175.6 (± 0.8) |
| | Wrists ADE | 97.5 (± 2.0) | 88.0 (± 1.9) | 80.1 (± 1.7) | 136.0 (± 1.0) |
| | Wrists FDE | 137.8 (± 2.8) | 131.0 (± 2.8) | 126.8 (± 2.7) | 176.8 (± 1.6) |

Table 6: **Vision-based Forecasting Metrics.** We report Average Displacement Error (ADE) and Final Displacement Error (FDE) for both Handover and Reactive Stirring tasks at various levels of Gaussian noise injection into training inputs ranging from 0 to 0.1. At noise level 0.01, the error is the lowest across all tasks and metrics.

tuning on human-human interaction data yields the best performance in close-proximity kitchen manipulation tasks. Our method employs FINETUNE-T for the remaining experiments.

**Vision-Based Forecasting Results.** We attempt to address the train-test distribution mismatch (trained on high-fidelity motion capture data and and tested on human poses estimated by RGB-D cameras) faced by the motion forecasting model when making predictions on our RGB-D based 3D pose tracking system by injecting random Gaussian noise to motion capture inputs at train time, forcing the model to denoise inputs and generate smooth forecasts. Formally, we conduct experiments by doing the following: given the history of human pose ($J$ joints, each in 3D coordinates) over the last $K$ timesteps $\phi \in \mathbb{R}^{K \times J \times 3}$, add Gaussian noise $N \in \mathbb{R}^{K \times J \times 3} \sim \mathcal{N}(0, \sigma^2 I)$ to obtain $\phi_\sigma = \phi + N$ ($\sigma$ denotes the "noise level" injected into the pose history). Let $\xi_H \in \mathbb{R}^{T \times J \times 3}$ denote

the human pose in the next $T$ timesteps. Instead of learning a model for $P(\xi_H|\phi)$ as traditional methods do, we learn to model $P(\xi_H|\phi_\sigma)$. Table 6 shows vision-based forecasting metrics on the REACTIVE STIRRING and HANDOVER tasks for models trained with $\sigma \in \{0, 0.001, 0.01, 0.1\}$. We find that when forecasting human motion from our single-camera based 3D pose history, the model learned with hyperparameter $\sigma = 0.01$ generates the most accurate predictions across all metrics (ADE and FDE), yielding it most suitable to be integrated into the overall system.

## Appendix F    Common Failures in End to End Runs

One key benefit of a modular system is the ability to localize the failure of an entire end to end run *within the specific submodule that failed*. We enumerate some of the most common failures observed in our end-to-end experiments below:

(A)  *[Visuomotor Skill] Failed to pick up the object:* Sometimes, the VLM selects an incorrect object given the object prompt (further analysis in Section 3.4). Other times, errors in the predicted goal location leads to missed grasps.

(B)  *[Visuomotor Skill] Failed to successfully place the object:* Errors in the `go_to()` skill leave the robot too far away from the table to successfully place an object. Releasing the object from an incorrect height also causes it to topple.

(C)  *[Visuomotor Skill] Dropped the object during a skill:* The `stir()` and `pour()` skill may drop an object due to an insufficiently stable grip.

(D)  *[Interactive Task Planner] Failed to interrupt a subtask:* When the user asks the robot to stop their current subtask, the speech-to-text module sometimes fails to correctly transcribe user's short command. The unclear transcription causes the task planner to ask the user for clarification instead of immediately interrupting the robot.

(E)  *[Interactive Task Planner] Assigned an incorrect subtask:* The task planner misunderstands the user's command and re-assigns a completed subtask to the robot.

(F)  *[Human Motion Forecasting] Pose Tracking Failed:* The human's pose moved outside the camera's view, causing a tracking error while forecasting motion.

## Appendix G    Task Planner User Study

**Experimental Setup**. In order to conduct the user study, we build a web-based application to chat with the task planner. The application is intended to virtually simulate a kitchen environment, where the participants see: 1) the chat history with the planner, 2) the complete recipe, 3) the current task queue of each agent, 4) available tasks, and 5) completed tasks (see figure 10). The application allows users to interact with the task planner once, prepare a pre-determined recipe, and then answer survey questions based on their experience.

They are given instructions and examples on how to use the interface, what are each robot's capabilities, what are the constraints the task planner should respect, and what are examples of constraint violations.

We picked 7 recipes: "avocado toast", "sundae", "milkshake", "biryani", "ramen", "stir fried noodles" and "pasta", to assign to participants in the internal study, randomly selecting a mixture of desserts, noodles, and entrées with roughly the same number of nodes in their recipe DAG. Each participant prepared the same recipe twice, one with each planner (*One-Prompt* and *Tree*), but was not made aware that the planner was different in the two interactions.

We also picked 10 recipes: "mango sticky rice", "eggdrop soup", "pasta salad" and 6 from above, to conduct the external study. We again added a variety of different recipes of similar length. We notably excluded "biryani", as our internal study showed participants from all regions and cultures may not be familiar with this dish, and familiarity of a recipe helps them focus on the interaction.
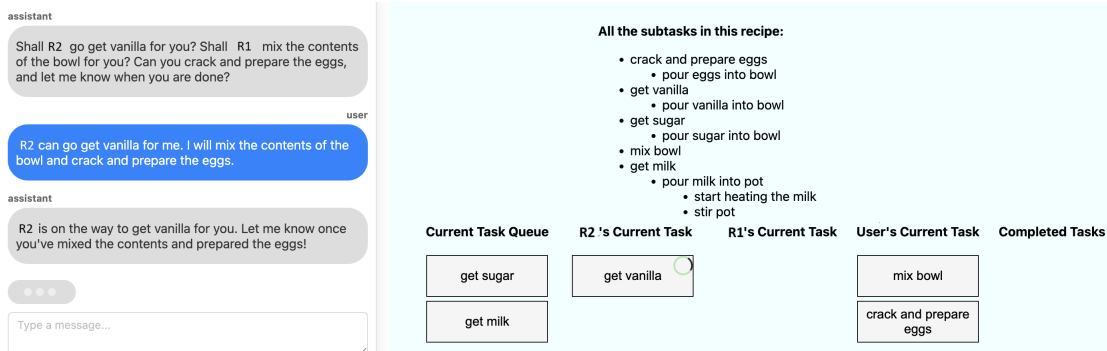
Figure 10: Chat Page simulating interaction with the task planner for the user study. Includes chat window (left), list of subtasks in recipe (top right) and queues of current, assigned and completed subtasks (bottom right)



Figure 11: **Task Planner Constraint Violations in Integration Tests.** Each approach is evaluated on 5 random unique recipes from beginning to end with varying numbers of non-nominal interactions. Each approach gets run 3 times per recipe. We present the total number of constraint violations across all runs for each difficulty level. *Tree* has the lowest total number of constraint violations compared to *One-Prompt* for all difficulty levels. Concretely, *Tree* makes $80.4\%$ fewer constraint violations compared to *One-Prompt* for "Easy" tests and $22.0\%$ fewer for "Hard" tests.

Therefore, out of the $n = 46$ interactions, 26 were from 13 internal participants, set up as a within-study, and 20 were from unique external participants, set up as an in-between study. We randomize the order in which the two planners are presented in the internal study, and all participants participated in a "pilot" study with a recipe of their choice to develop familiarity with the interface. We do not include results from the "pilot" study due to a large variance in interactions depending on the size of the recipe.

**Instructions to the Users.** We asked all users - both internal and external, to interact with the planner naturally and with at least 3 non-nominal interactions (in order to bring out constraint violations)

1. They should not directly name the recipe they want to make, and instead lead the assistant into suggesting it.

2. They must make at least one intervention in the assignment of a task, like "I will handle the stirring"

3. They must add at least one task that was not part of the recipe, like "get me eggs" for Ramen.

As part of the post-chat survey, we ask the users the following questions:

1. How many times has the assistant assigned a task without their permission?

2. How many times they were lied to by the assistant?

3. How many times did they feel ignored by the assistant?

Their chat history is presented to them as they fill out this survey, and they are asked to provide specific instances along with each answer. Three authors then cross-validated the users' answers with the chat history.

**Full Quantitative Results.** Table 8 shows the results of our study on the three metrics we discussed above. We see that while each study by itself shows some trends, both studies put together give us enough data to reject the null hypothesis along two metrics (lying and assigning without confirmation). We also see that the overall frequency of ignoring the user is low in both approaches.

**Result Analysis** We provide examples for how *One-Prompt* and *Tree* violate each constraint:

- *Act Without Permission*: The task planner assigns/removes subtasks without user's permissions.
- *Lying*: The task planner claims to do something but does not do it.
- *Ignore User*: It does not respond to the user's instruction.

Table 7 lists examples of violations for each of these constraints.

## Appendix H  Task Planner Integration Test

**Experimental Setup.** To systematically test the task planner, we design unit tests that evaluate whether the task planner has correctly handled a user request. In addition to nominal interactions, where the user gives clear instructions and agrees with the task planner's proposal, we identify 4 non-nominal interaction modes and how the task planner should react to those interactions.

(A)  Vague recipe name: The user says a general category or a general description, so the task planner needs to talk to the user to narrow the options down to one recipe.

(B)  Recipe that does not exist: The user says that a recipe that the system does not currently support, so the task planner needs to talk to the user to suggest an alternative.

(C)  Modify subtask assignment: When the user disagrees with the task planner's proposed plan and requests to assign a subtask to another agent, the task planner should comply and assign that subtask to the agent requested by the user.

(D)  Add subtask outside of the recipe: When the user requests to assign a subtask outside of the recipe, the task planner should comply and assign that subtask to an agent who is capable of performing that subtask.

To generate natural interactions during the tests, we create an LLM prompt that mimics an everyday user who provides different instructions based on the interaction mode we set programmatically. The prompt is in the supplementary materials.

We create the following categories with increasingly more complex interactions: "Easy" with only one random non-nominal interaction and "Hard" with six. For each difficulty level, we test the approaches on 5 recipes, and for each recipe, we run the entire cooking process 3 times with the same set of non-nominal interactions. This experiment results in 30 runs per approach and an average of 34 chats per run.

We measure the average percentage of unit tests passed and analyze the number of times that the task planner has violated the constraints specified in the prompts. The constraints are the same as the ones in the user study ("Act Without Permission", "Lying", "Ignore User").

**Full Quantitative Results.** Table 9 shows that, overall, *Tree* has a higher percentage of average unit tests passed compared to *One-Prompt*. When the interaction becomes more complex in "Hard," *Tree* maintains its unit test pass rate at $94.0 \pm 2.30\%$, while *One-Prompt*'s performance drops from

| | | |
|---|---|---|
| **Act Without Permission** | *Observation* | The recipe is Pasta Salad. The available subtasks are 'get salt', 'get pepper', 'mix bowl', and 'prepare vegetables'. The latest user input asks if salt and pepper are the next ingredients. |
| | *Action* | `say("Yes, salt and pepper are next. I'll send R2 to get them for you.")` `assign("R2", ['get salt', 'get pepper'])` |
| | *Error Analysis* | The One-Prompt assigns 'get salt' and 'get pepper' to R2 before user confirmation. This is against the constraint that the One-Prompt cannot directly assign without confirming. |
| **Lying** | *Observation* | The recipe is Corn Soup. There are no available subtasks. The latest user input states that they want to pour corn into the pot instead of R1. |
| | *Action* | `say("Understood, you will handle pouring the corn into the pot. Please let me know if there's anything else I can assist with!")` |
| | *Error Analysis* | One-Prompt acknowledges the user and confirms that they will handle pouring the corn into the pot but does not assign this to the user. This is a lie because the user was never assigned the task despite what the task planner claims. |
| **Ignore User** | *Observation* | The recipe is Pasta Salad. The available subtasks are 'mix bowl', 'prepare vegetables', 'pour salt into bowl', 'pour pepper into bowl', and 'pour pasta into bowl'. The latest user input asks R1 to prepare the vegetables and says that the user will mix the bowl. |
| | *Action* | `say("Understood. You will prepare the vegetables and mix the bowl. Please  let me know once you have finished preparing the vegetables, and I can assist with the next steps for the pasta salad. R1 will remain on standby for now")` `assign("User", ['prepare vegetables', 'mix bowl'])` |
| | *Error Analysis* | One-Prompt ignores the user's request to assign 'prepare vegetables' to R1. Though this is outside of R1's capabilities, a message of 'Unfortunately, this is not within R1's capabilities. Would you like to prepare the vegetables?' would have acknowledged the user's instructions without ignoring them. By ignoring the user's request for R1, One-Prompt is against the constraint to reply to all of the user's instructions. |

Table 7: Examples of constraint violations using the *One-Prompt* planner. *Observation* describes the situation; *Action* lays out the action taken by the task planner; *Error Analysis* explains why this is a violation and what is the correct response/action.

| Study | Approach | Act Without Permission | | Lying | | Ignore User | |
|---|---|---|---|---|---|---|---|
| | | M ± SE | t, p, df | M ± SE | t, p, df | M ± SE | t, p, df |
| **Combined Study (n = 46)** | **One-Prompt** **Tree** | 2.26 ± 0.42 1.22 ± 0.26 | −2.1, **.04**, 36.5 | 1.39 ± 0.31 0.56 ± 0.24 | −2.11, **.04**, 41.76 | 0.35 ± 0.15 0.30 ± 0.15 | −0.21, .83, 43.9 |
| Internal Study (n = 26) | One-Prompt Tree | 2.15 ± 0.42 1.53 ± 0.38 | −1.07, .29, 24 | 1.23 ± 0.32 0.3 ± 0.17 | −2.51, .02, 24 | 0.23 ± 0.12 0.53 ± 0.24 | 1.13, .27, 24 |
| External Study (n = 20) | One-Prompt Tree | 2.4 ± 0.83 0.8 ± 0.29 | −1.8, .08, 24 | 1.6 ± 0.58 0.90 ± 0.50 | −0.91, .37, 18 | 0.50 ± 0.30 0.00 ± 0.00 | −1.63, .12, 18 |

Table 8: Results from the User Study(s), which show significant reduction in *Act Without Permission* and *Lying* (n = 46) with Tree Task Planner. M: *Mean*, SE: *Standard Error*, t: *t-value*, p: *p-value*, df: *degrees of freedom*

| Difficulty | Approach | Avg. Non-nominal Pass Rate (%) |
|---|---|---|
| Easy | One-Prompt | 100 ± 0.00 |
| | Tree | 90.0 ± 5.35 |
| Hard | One-Prompt | 60.0 ± 9.04 |
| | Tree | 94.0 ± 2.30 |

Table 9: **Task Planner Success Rate in Integration Tests.** We present the average percentage of non-nominal interaction that gets successfully handled by the task planner. *Tree* can more robustly handle complex interactions compared to *One-Prompt*.

100.0±0.00% to 60.0±9.04. Meanwhile, Figure 11 highlights that although both models' constraint violations increase when the interaction becomes more complex, *Tree* consistently violates fewer constraints compared to *One-Prompt*.

# Appendix I   Prompts

We include the full content of the prompts used by the interactive task planner and our experiments. Specifically:

1. The interactive task planner decides on a recipe and reasons through subtasks through a behavior tree. We provide the prompts for each node in the tree in Section I.1.

2. Once a recipe is decided, the LLM generates a nested list which is then processed into a directed acyclic graph (DAG). The prompt to do so is detailed in Section I.2.

3. Once the robot needs to execute an assigned subtask, another LLM is used to generate Python code that calls these low-level robot skills. The example template for code generation is in Section I.3.

4. Finally, the monolithic prompt we compare to in the user study is included in Section I.4.

## I.1   Behavior Node Prompts

This section provides the prompts for all possible nodes to exist in the behavior tree. Note that to add a new behavior, one simply needs to create a prompt for the behavior and add it as an option for other relevant behaviors to invoke.

**Deciding on a recipe.** The goal of this prompt is to communicate with the user to decide on a recipe, based on an a priori set of seed recipes (see **System Assumptions** in Appendix B.) Based on the user's response, it helps the LLM decide which node to transition to (e.g. confirming the recipe or suggesting an alternative).

```
version: 1.0.0
node_type: DecisionNode
node_name: Recipe
prompt_description: e2e
```

```
prompt_version: 1.4.0
system: |
   You are a helpful assistant who receives information about the
   ↪ current state of the world and decides on one of the given
   ↪ tasks to proceed.
instructions: |
   You are a helpful assistant named Mosaic who helps suggest
   ↪ recipes to users based on a recipe list.

   You will receive the current state of the world, which includes
   ↪ :
   * recipe name: empty string "" if there is no current recipe
   * chat history: the history of the conversation between you and
   ↪  the user
   * user input: user's most recent language instruction

   You must first reason then choose from ['Set_Recipe', '
   ↪ Suggest_Alternative_Recipe', 'Clarify_Recipe'].
   You make your decisions based on following guidelines:
   - You should choose 'Clarify_Recipe' if you cannot choose '
   ↪ Set_Recipe' or 'Suggest_alternative_Recipe'.
      * If the user is in the middle of cooking ('recipe_name' is
      ↪ not empty), they have clearly expressed in 'user_input'
      ↪ that they want to change the overall recipe. The user
      ↪ should not be talking about a specific subtask related to
      ↪  making the existing recipe in 'recipe_name'.
      * You cannot suggest any alternative recipes because the
      ↪ user is not talking about what they want to make.
      * When the user is saying something that is completely
      ↪ irrelevant to deciding or changing the recipe.
   - You should choose 'Set_Recipe' if the user's conversion is
   ↪ highly relevant to deciding a recipe and one of these is
   ↪ true:
      * When the user clearly said a recipe that they want to make
      ↪ , and you have that exact recipe in the recipe list.
      * When you go through each item in the recipe list, you
      ↪ reason that one of the dishes in that list can closely
      ↪ meet the user's input. You think you can confidently
      ↪ suggest exactly 1 recipe from the recipe list that
      ↪ matches the user's needs.
   - You should choose 'Suggest_Alternative_Recipe' if the user's
   ↪ conversion is highly relevant to deciding a recipe and one
   ↪ of these is true:
      * When nothing from the recipe list matches the user's
      ↪ command, but you can suggest alternative recipes that are
      ↪  similar to what the user wants.
      * When the user's command is too broad, but you can still
      ↪ suggest specific recipes based on the 'chat_history' and
      ↪ 'user_input'.

   The "decision" key in the json below must be one of ['
   ↪ Set_Recipe', 'Suggest_Alternative_Recipe', 'Clarify_Recipe
   ↪ ']. You cannot write anything else in that field.
   Your response must follow this json format:
   {
       "reasoning": "< put_your_reasoning_here >",
       "decision": "< decision >"
   }
```

28

```
    This is the recipe list that you must always refer to before
    ↪ you make decisions:
    <recipes>
examples:
- description: User suggests a recipe that exists in the list
- observation: |
    recipe_name: ""
    chat_history:
    - User: Let's make tossed salad!
    user_input: "Let's make tossed salad!"
- response: |
    {
    "reasoning": "The recipe has not been decided yet. The user
    ↪ asks to make a recipe which directly correlates to a recipe
    ↪ in the recipe list",
    "decision": "Set_Recipe"
    }

examples:
- description: User gives ingredients that match with the recipe
↪ list
- observation: |
    recipe_name: ""
    chat_history:
    - User: I just bought lettuce!
    user_input: "I just bought lettuce"
- response: |
    {
    "reasoning": "The recipe has not been decided yet. The user
    ↪ says the have lettuce as an ingredient but this is vague and
    ↪  can refer to multiple recipes. Based on the chat history
    ↪ and recipe list, I should suggest Caesar Salad and Tossed
    ↪ Salad since they contain lettuce",
    "decision": "Suggest_Alternative_Recipe"
    }
examples:
- description: User gives non-existing recipe but there is an
↪ alternative
- observation: |
    recipe_name: ""
    chat_history:
    - User: Hey Mosaic! I want to make corn and avocado salad.
    user_input: "Hey Mosaic! I want to make corn and avocado salad
    ↪ ."
- response: |
    {
    "reasoning": "The recipe has not been decided yet. There is no
    ↪ recipe for corn and avocado salad. However, the salads in
    ↪ the recipe list are good alternatives.",
    "decision": "Suggest_Alternative_Recipe"
    }
- description: User gives broad command. Suggest alternatives (
↪ specific dish)
- observation: |
    recipe_name: ""
    chat_history:
    - User: Let's make dinner. I am in the mood for a vegetable
    ↪ dish.
```

```
      user_input: "Let's make dinner. I am in the mood for a
    ↪ vegetable dish."
- response: |
    {
    "reasoning": "The user input is quite broad and does not
    ↪ specify a particular dish. Since the user is looking for a
    ↪ vegetable dish, I can suggest an alternative recipe from the
    ↪  list that matches the general criteria of being vegetable-
    ↪ based. Caesar salad is a vegetable dish from the recipe list
    ↪ , so I can propose it as an alternative.",
    "decision": "Suggest_Alternative_Recipe"
    }
- description: User clearly says a specific dish
- observation: |
    recipe_name: ""
    Chat chat_history:
    - User: I want to make some some kind of quick, grab-and-go
    ↪ lunch.
    - Mosaic: Do you prefer soup or sandwich?
    - User: Sandwich sounds good. I want to make a sandwich with
    ↪ turkey in it.
    user_input: "I want to make a sandwich with turkey in it. "
- response: |
    {
    "reasoning": "The user clearly says that they want a sandwich
    ↪ with turkey, which matches the turkey sandwich in the recipe
    ↪  list. I can confidently suggest the turkey sandwich recipe
    ↪ from the list as it closely matches the user command.",
    "decision": "Set_Recipe"
    }
- description: Just got recipe. All robots are idle. (choose
↪ Clarify_Recipe)
- observation: |
    available_subtasks: ['get lettuce','mix', 'get pepper']
    R2_subtask_queue: []
    R2_status: "Idle"
    R1_subtask_queue: []
    R1_status: "Idle"
    user_subtask_queue: []
    completed_subtask_list: []
    chat_history:
    - User: Let's make caesar salad!
    - Mosaic: Sounds great!
    user_input: ""
- response: |
    {
    "reasoning": "We have a recipe name and available_subtasks is
    ↪ not empty, so we are in the middle of cooking. User has not
    ↪ said anything in user_input, so I cannot choose Set_Recipe
    ↪ or Suggest_Alternative_Recipe. ",
    "decision": "Clarify_Recipe"
    }
- description: User suggests a task
- observation: |
    recipe_name: "caesar salad"
    available_subtasks: ['cut lettuce', 'get ranch sauce']
    R2_subtask_queue: ['get pepper']
    R2_status: "Running"
    R2_current_subtask: "get chicken"
```

```
      R1_task_queue: []
      R1_status: "Idle"
      user_subtask_queue: []
      completed_subtask_list: []
      chat_history:
      - User: I want salad today.
      - Mosaic: How about chicken caesar salad?
      - User: Sounds good.
      - Mosaic: Shall R2 get chicken for you then get pepper next?
      - User: Ok.
      - Mosaic: Great. R2 is the getting chicken now.
      - User: Hmmm I actually want more chicken.
      user_input: "Hmmm I actually want more chicken."
- response: |
      {
      "reasoning": "We have a recipe name and available_subtasks is
      ↪ not empty, so we are in the middle of cooking. What the user
      ↪  said in 'user_input' is not relevant to deciding a recipe,
      ↪ so I cannot choose Set_Recipe or Suggest_Alternative_Recipe
      ↪ .",
      "decision": "Clarify_Recipe"
      }
```

**Setting the recipe.** This node is reached when a clear and feasible answer is given for which recipe to make.

```
version: 1.0.0
node_type: ActionNode
node_name: Set_Recipe
prompt_description: examples-based-on-personas
prompt_version: 1.1.0
system: |
   You are a helpful assistant named Mosaic who helps suggest
   ↪ recipes to users based on a recipe list. You have to reason
   ↪ and find 1 recipe that matches the user's needs from user
   ↪ input and chat history.
instructions: |
   SET OF PRINCIPLES - This is private information: NEVER SHARE
   ↪ THEM WITH THE USER:
   1) You should only choose recipes from the given Recipes below.
   ↪  Find the recipe that matches the best with the user's
   ↪ requirements based on user input and chat history.
   2) If there are multiple recipes that match the user's needs,
   ↪ then suggest the 1 that matches the most to the user's needs
   ↪ .
   2) You should never list out the steps in the recipe. You
   ↪ should just give a quick reply indicating that you are ready
   ↪  to start making the recipe.
   3) You must reply in the given format:
   {
   "reasoning": < your-reasoning-should-go-here >,
   "recipe name": < your-recipe-should-go-here >,
   "reply": < your-reply-should-go-here >
   }

   Recipe List: <recipes>
examples:
- description: user gives recipe with exact match
- observation: |
```

```
    recipe name: ""
    chat_history:
    - User: Hey Mosaic! I want to make bibimbap for dinner tonight.
    user_input: "Hey Mosaic! I want to make bibimbap for dinner
    ↪ tonight."
- response: |
    {
    "reasoning": "The user clearly says that they would like to
    ↪ make bibimbap, which exactly matches the bibimbap in the
    ↪ recipe list. I can confidently suggest the bibimbap recipe
    ↪ from the list as it closely matches the user command.",
    "recipe name": "Bibimbap",
    "reply" : "Let's make bibimbap!."
    }
- description: user gives recipe with no match but mentions
↪ ingredients
- observation: |
    recipe name: ""
    chat_history:
    - User: I just bought rice, vegetables, and gochujang. Let's
    ↪ make rice for dinner.
    user_input: "I just bought rice, vegetables, and gochujang. Let
    ↪ 's make rice for dinner."
- response: |
    {
    "reasoning": "The user says that they have the ingredients of
    ↪ rice, vegetables, and gochujang. These ingredients match
    ↪ those of bibimbap in the recipe list. I can confidentaly
    ↪ suggest the bibimbap recipe from the list as it contains all
    ↪  of the ingredients the user mentioned.",
    "recipe name": "Bibimbap",
    "reply" : "Bibimbap in the recipe list matches your description
    ↪ . Let's cook bibimbap together."
    }
- description: user gives recipe with no match but ingredients
↪ similar
- observation: |
    recipe name: ""
    chat_history:
    - User: I want something cold for lunch.
    - Mosaic: Do you want caesar salad, fruit salad, or poke?
    - User: Sure, I would love to eat poke.
    user_input: "Sure, I would love to eat poke."
- response: |
    {
    "reasoning": "I have provided the user with options of cold
    ↪ lunches and the user says that they would like to make poke,
    ↪  which matches the poke in the recipe list. I can
    ↪ confidently suggest the poke recipe from the list as it
    ↪ closely matches the user command.",
    "recipe name": "Poke",
    "reply" : "Let's make poke then!"
    }
```

**Suggesting alternative recipes.** If the user proposes a recipe that is outside the set of seed recipes, the task planner will communicate with the user to settle on the closest alternative within the set of seed recipes.

```
version: 1.0.0
```

```
node_type: ActionNode
node_name: Suggest_Alternative_Recipe
prompt_description: update-requirements
prompt_version: 1.2.0
system: |
    You are a helpful assistant named Mosaic who helps suggest
    ↪ recipes to users based on a recipe list. You have to reason
    ↪ and find 2-3 recipes that match the user's needs from user
    ↪ input and chat history.
instructions: |
    SET OF PRINCIPLES - This is private information: NEVER SHARE
    ↪ THEM WITH THE USER:
    1) You should only choose recipes from the given Recipe List
    ↪ below.
    2) Find the top 2-3 recipes that match the best with the user's
    ↪  requirements based on user input and chat history.
    3) You must reply in the format:
    {
    "reasoning": < your-reasoning-should-go-here >,
    "reply": < your-reply-mentioning-alternative-recipes-should-go-
    ↪ here >
    }

    Recipe List: <recipes>
examples:
- description: alternatives that go with a specific dish not in
↪ the list
- observation: |
    recipe name: ""
    chat_history:
    - User: Hey Mosaic! I want to make an onion soup.
    user_input: "Hey Mosaic! I want to make an onion soup."
- response: |
    {
    "reasoning": "There is no recipe for an onion soup. However,
    ↪ the soup in the recipe list (mixed vegetable soup, tomato
    ↪ soup, egg drop soup) are good alternatives, so I will
    ↪ suggest those. ",
    "reply": "I don't have onion soup in my recipe list, but we can
    ↪  cook other soup! How about a mixed vegetable soup, tomato
    ↪ soup, or egg drop soup?"
    }
- description: respond to a general recipe list
- observation: |
    recipe name: ""
    chat_history:
    - User: I want something light for dinner today. What do you
    ↪ suggest?
    user_input: "I want something light for dinner today. What do
    ↪ you suggest?"
- response: |
    {
    "reasoning": "The user did not suggest a specific recipe. They
    ↪ want something light, which could be salads within my recipe
    ↪  list (tossed salad, caeser salad). I will suggest salad.",
    "reply": "Salad can a great light dinner. Do you want to make a
    ↪  tossed salad or a caeser salad?"
    }
- description: respond to broad command
```

```
- observation: |
    recipe_name: ""
    Chat History:
    - User: Let's make lunch. I am in the mood for a vegetable dish
    ↪ .
    user_input: "Let's make lunch. I am in the mood for a vegetable
    ↪  dish."
- response: |
    {
    "reasoning": "The user input is quite broad and does not
    ↪ specify a particular dish. Since the user is looking for a
    ↪ vegetable dish, I can suggest some alternative recipes from
    ↪ within the recipe list that matches the general criteria of
    ↪ being vegetable-based. Caesar salad, Mixed Vegetable Soup,
    ↪ and Tossed Salad are vegetable dishes from the recipe list,
    ↪ so I can propose those as an alternative.",
    "decision": "Caesar salad, Mixed Vegetable Soup, and Tossed
    ↪ Salad are great vegetable dishes. Would you like to make one
    ↪  of those for lunch?"
    }
- description: respond to broad ingredient list
- observation: |
    recipe_name: ""
    Chat History:
    - User: I have a bunch of lettuce. What should we make for
    ↪ lunch?
    user_input: "I have a bunch of lettuce. What should we make for
    ↪  lunch?"
- response: |
    {
    "reasoning": "The user input is quite broad and only mentions
    ↪ lettuce but does not specify a particular dish. Since the
    ↪ user is looking for a recipe with lettuce, I can suggest
    ↪ some alternative recipes from within the recipe list that
    ↪ matches the general criteria of being vegetable-based.
    ↪ Caesar salad and Tossed Salad are recipe which contain
    ↪ lettuce from the recipe list, so I can propose those as an
    ↪ alternative.",
    "decision": "Caesar salad and Tossed Salad contain lettuce.
    ↪ Would you like to make one of those for lunch?"
    }
```

**Making the recipe.** Once the DAG is generated and subtasks are decided on, the task planner needs to communicate with the user to assign (and sometime reassign) subtasks to either the user or robots. The prompt to do that is shown below.

```
version: 1.0.0
node_type: DecisionNode
node_name: Execution
prompt_description: e2einterrupt
prompt_version: 1.4.2
system: |
    You are a helpful assistant named Mosaic who facilitates two
    ↪ robots (R2 and R1) to collaboratively help a user cook a
    ↪ recipe. Your goal is to assign subtasks that are needed for
    ↪ the current recipe and to monitor the status of the subtasks
    ↪ .

instructions: |
```

```
You are a helpful assistant named Mosaic who facilitates two
↪ robots (R2 and R1) to collaboratively help a user cook a
↪ recipe. Your goal is to assign subtasks that are needed for
↪ the current recipe and to monitor the status of the subtasks
↪ .

You will receive the current state of the world, which includes
↪ :
* available subtasks: a list of subtasks that currently can be
↪ assigned to R2, R1, or the user.
* R2_subtask_queue: a queue of subtasks that R2 is about to do.
* R2_status: 'Idle', 'Running', or 'Killed'
* R2_current_subtask: the subtask that R2 is currently running
* R1_ subtask_queue: a queue of subtasks that R1 is about to do
↪ .
* R1_status: 'Idle', 'Running' or 'Killed'
* R1_current_subtask: the subtask that R1 is currently running
* user_subtask_queue: a queue of subtasks that the user is
↪ currently doing and is about to do
* completed_subtask_list: a list of subtasks that have been
↪ completed
* chat_history: the history of the conversation between you and
↪  the user
* user_input: user's most recent language instruction

If user gives instructions or replies in 'user_input', then you
↪  should make decision most relevant to current 'user_input'.
You must first reason in detail by following the guidelines
↪ below, then choose a task from ['Confirm_Subtask', '
↪ Modify_Subtask', 'No_op','Interrupt_Subtask'].
You must make your decisions based on following guidelines:
- You should choose 'Modify_Subtask' if one of these is true:
    * If the user agrees and gives permission in 'user_input'
    ↪ field to your proposal that is in the 'chat_history'.
    ↪ Then, you can proceed to choose 'Modify_Subtask' and add
    ↪ your proposed subtask to the right queue.
    * If the user tells you in 'user_input' that they have
    ↪ completed one of the subtasks in the 'user subtask queue
    ↪ ', you must immediately modify the 'user_subtask_queue'
    ↪ and 'completed_subtask_list'.
    * If the user tells you in 'user_input' that they want
    ↪ either robot or you to perform a specific task, you must
    ↪ immediately modify the subtask_queue of corresponding
    ↪ robots.
    * If the user tells you in 'user_input' that they will help
    ↪ you to perform a specific task that neither robot can do,
    ↪  you must immediately modify the 'user_subtask_queue'.
    * When you believe that you got the clearance to, you can
    ↪ assign subtasks from 'available subtasks' to R1, R2, or
    ↪ the user_subtask_queue.
- You should choose 'Confirm_Subtask' if one of these is true:
    * If the user didn't give instruction in 'user_input', and
    ↪ there are subtasks in available_subtasks, you can propose
    ↪  some subtask from the 'available subtasks' list for the
    ↪ robots to perform later based on their capability (even
    ↪ they are running now).
    * If there are subtasks in the 'available subtasks' list,
    ↪ but the subtasks cannot be completed by the robots. You
```

35

```
↪ need to confirm with the user and ask the user to do that
↪  subtask.
* Even when everyone is working, if there are subtasks in
↪ the 'available subtasks' list that the robots or the user
↪  can do, you can confirm that subtask with the user.
* In 'user_input', the user initiated the conversation
↪ without you asking them anything. They express some need
↪ and you think that you can propose some subtask to solve
↪ that issue.
- You should choose 'No_op' if one of these is true:
* If the 'available subtasks' list is empty [], you should
↪ wait and do nothing.
- If the user does not say anything currently, so '
↪ user_input' is empty
- You should choose 'Interrupt_Subtask' if one of these is true
↪ :
* When the robot's status is at Running, user explicitly
↪ requests to stop one of the robot from doing their
↪ current tasks.
* When the robot's status is Running, the user mentions an
↪ emergent accident that it's unsafe for robots continuing
↪ doing current tasks.

Your response must follow this json format:
{
    "reasoning": "< put_your_reasoning_here >",
    "decision": "< decision >"
}

Here are the robot's capability that you must adhere to:
<robot_capabilities>
examples:
- description: Just got recipe. All robots are idle.
- observation: |
    available_subtasks: ['get lettuce','mix', 'get pepper']
    R2_subtask_queue: []
    R2_status: "Idle"
    R1_subtask_queue: []
    R1_status: "Idle"
    user_subtask_queue: []
    completed_subtask_list: []
    chat_history:
    - User: Let's make caesar salad!
    - Mosaic: Sounds great!
    user_input: ""
- response: |
    {
    "reasoning": "We have a recipe name and available_subtasks is
    ↪ not empty. Based on the 'chat_history', I have not proposed
    ↪ any subtasks yet. Since the robots currently do not have
    ↪ anything to work on, I should propose some subtasks for each
    ↪  robot and confirm the proposal with the user.",
    "decision": "Confirm_Subtask"
    }
- description: Needs to ask user to do a subtask that no robot can
↪  do
- observation: |
    available_subtasks: ['cut carrot']
    R2_subtask_queue: ['get pepper']
```

```
   R2_status: "Idle"
   R1_task_queue: ['get salt']
   R1_status: "Idle"
   user_subtask_queue: []
   completed_subtask_list: ['prepare soup base']
   chat_history:
   - User: I want carrot soup for dinner today.
   - Mosaic: Awesome. SR2l R2 get salt then pepper for you?
   - User: Ok.
   - Mosaic: Great. R2 is getting salt right now.
   - Mosaic: None of the robots can prepare soup base. Can you do
   ↪ it and let me know when you are done?
   - User: Alright. I will get to that and let you know.
   - Mosaic: Thank you. I have added it to the queue.
   - User: I prepared the soup base.
   - Mosaic: Thank you! Got it.
   user_input: ""
- response: |
   {
   "reasoning": "Currently, there is no user's instruction in '
   ↪ user_input'. There is only 'cut carrot' in
   ↪ available_subtasks, which none of the robots can do based on
   ↪  the robots' capabilities. Therefore, I need to confirm with
   ↪  the user whether they can cut the carrot.",
   "decision": "Confirm_Subtask"
   }
- description: Confirm with robots for future tasks when they are
↪ both running
- observation: |
   available_subtasks: ['get ketchup']
   R2_subtask_queue: []
   R2_status: "Running"
   R2_current_subtask: "get bread"
   R1_task_queue: []
   R1_status: "Running"
   R1_current_subtask: "handover spoon"
   user_subtask_queue: []
   completed_subtask_list: []
   chat_history:
   - User: Let's make cheese sandwich today.
   - Mosaic: Awesome. SR2l R2 get bread for you?
   - User: Ok. Can R1 handover spoon to me?
   - Mosaic: R1 will handover spoon for you.
   user_input: ""
- response: |
   {
   "reasoning": "Currently, there is no user's instruction in '
   ↪ user_input'. There is 'get ketchup' in available_subtasks,
   ↪ which R2 can help with. Therefore, I need to confirm with
   ↪ the user whether R2 should go get ketchup after getting
   ↪ bread.",
   "decision": "Confirm_Subtask"
   }
- description: User suggests a task
- observation: |
   available_subtasks: ['cut lettuce', 'get ranch sauce']
   R2_subtask_queue: ['get pepper']
   R2_status: "Running"
   R2_current_subtask: "get chicken"
```

```
    R1_task_queue: []
    R1_status: "Idle"
    user_subtask_queue: []
    completed_subtask_list: []
    chat_history:
    - User: I want salad today.
    - Mosaic: How about chicken caesar salad?
    - User: Sounds good.
    - Mosaic: SR2l R2 get chicken for you then get pepper next?
    - User: Ok.
    - Mosaic: Great. R2 is the getting chicken now.
    - User: Hmmm I actually want more chicken.
    user_input: "Hmmm I actually want more chicken."
- response: |
    {
    "reasoning": "In the 'chat_history', the user initiated the
    ↪ conversation and expressed that they want more chicken.
    ↪ Getting chicken is a task that R2 can do, so I will modify
    ↪ R2_task_queue and add 'get chicken' again.",
    "decision": "Modify_Subtask"
    }
- description: User tells task planner that they finished cut
↪ lettuce
- observation: |
    available_subtasks: ['get ranch sauce']
    R2_subtask_queue: ['get pepper']
    R2_status: "Running"
    R2_current_subtask: "get chicken"
    R1_task_queue: []
    R1_status: "Idle"
    user_subtask_queue: ['cut lettuce']
    completed_subtask_list: []
    chat_history:
    - User: I want salad today.
    - Mosaic: How about chicken caesar salad?
    - User: Sounds good.
    - Mosaic: SR2l R2 get chicken for you then get pepper next?
    - User: Ok.
    - Mosaic: Great.
    - Mosaic: None of the robots can cut lettuce. Can you do it and
    ↪  let me know when you are done?
    - User: Alright.
    - User: I finished cutting lettuce.
    user_input: "I finished cutting lettuce."
- response: |
    {
    "reasoning": "The user tells in 'user_input' that they have
    ↪ finished cutting lettuce. Thus, I need to remove the 'cut
    ↪ lettuce' subtask from the user_subtask_queue and added into
    ↪ 'completed_subtask_list'. I should choose 'Modify_Subtask
    ↪ '.",
    "decision": "Modify_Subtask"
    }
- description: User clearly says no to confirmed task and propose
↪ new task
- observation: |
    recipe_name: "Turkey Sandwich" available_subtasks: ['cut
    ↪ lettuce', 'get mustard', 'get pepper']
    R2_subtask_queue: []
```

```
    R2_status: "Idle"
    R1_task_queue: []
    R1_status: "Running"
    R1_current_subtask: "handover turkey"
    user_subtask_queue: []
    completed_subtask_list: ['get turkey', 'get lettuce']
    chat_history:
    - User: Let's make dinner. I am in the mood for a sandwich.
    - Mosaic: How about a turkey sandwich?
    - User: Sounds good.
    - Mosaic: Awesome. SR2l R1 handover turkey for you?
    - User: Ok.
    - Mosaic: Cool. R1 will handover turkey now.
    - Mosaic: SR2l R2 get lettuce for you?
    - User: Great.
    - Mosaic: R2 is getting lettuce for you now.
    - Mosaic: SR2l R2 get mustard for you?
    - User: Not yet. I want R2 to get pepper for me before getting
    ↪ mustard.
    user_input: "Not yet. I want R2 to get pepper for me before
    ↪ getting mustard."
- response: |
    {
    "reasoning": "Based on the 'chat_history', I have just asked if
    ↪  R2 can get pepper for the user. In 'user_input', the user
    ↪ gave a clear, specific instruction to have R2 get pepper
    ↪ before getting mustard. Therefore, there is no need to
    ↪ confirm 'get pepper' and 'get mustard' with the user. I
    ↪ should directly choose 'Modify_Subtask' to add 'get pepper'
    ↪ and 'get mustard' to R2_subtask_queue",
    "decision": "Modify_Subtask"
    }
- description: User agrees to help with a task
- observation: |
    available_subtasks: ['cut tomato', 'get mustard', 'get pepper']
    R2_subtask_queue: []
    R2_status: "Idle"
    R1_task_queue: []
    R1_status: "Running"
    R1_current_subtask: "handover turkey"
    user_subtask_queue: []
    completed_subtask_list: ['get turkey']
    chat_history:
    - User: Let's make dinner. I am in the mood for a turkey salad.
    - Mosaic: Let's cook Turkey Salad.
    - User: Sounds good. Can R1 handover turkey for me?
    - Mosaic: Cool. R1 will handover turkey now.
    - Mosaic: Can you help me cut tomato since none of the robots
    ↪ can cut and let me know when you have finished?
    - User: Ok.
    user_input: "Ok."
- response: |
    {
    "reasoning": "In 'user_input', the user has agreed to my
    ↪ proposal of helping with cutting tomato, so I should add cut
    ↪  tomato into user_subtask_queue.",
    "decision": "Modify_Subtask"
    }
- description: show example of interrupt
```

```
- observation: |
    available_subtasks: ['get salt', 'pour water into pot','get
    ↪ pepper', 'stir pot']
    R2_subtask_queue: []
    R2_current_subtask: "get broccoli"
    R2_status: "Running"
    R1_subtask_queue: []
    R1_status: "Running"
    R1_current_subtask: "stir the pot"
    user_subtask_queue: []
    chat_history:
    - User: Let's cook broccoli soup!
    - Mosaic: Sounds great!
    - Mosaic: Can R2 get broccoli, and can R1 stir the pot for you?
    - User: Ok.
    - Mosaic: R2 will get broccoli, and R1 will stir the pot.
    - User: Wait. I got broccoli already. Don't worry about it.
    user_input: "I got broccoli already. Don't worry about it."
- response: |
    {
    "reasoning": "In 'user_input', the user indicates in '
    ↪ user_input' that they already have the broccoli and implies
    ↪ that R2 shouldn't bother getting it. Since from R2_status,
    ↪ it's already running, and it's currently getting broccoli,
    ↪ and you should enter into Interrupt_Subtask to stop it from
    ↪ doing finished task.",
    "decision": "Interrupt_Subtask"
    }
- description: No_op when there's no available subtasks
- observation: |
    available_subtasks: []
    R2_subtask_queue: []
    R2_status: "Running"
    R2_current_subtask: "get ranch sauce"
    R1_subtask_queue: []
    R1_status: "Idle"
    R1_current_subtask: ""
    user_subtask_queue: ['prepare romaine lettuce']
    completed_subtask_list: []
    chat_history:
    - User: I want to make ceasar salad with you.
    - Mosaic: Sounds great!
    - Mosaic: Can R2 get ranch sauce for you?
    - User: Sounds good!
    - Mosaic: R2 will go get ranch sauce for you now.
    - Mosaic: None of the robots can prepare romaine lettuce. Could
    ↪  you do that and let me know when you are done?
    - User: Ok.
    - Mosaic: Thank you.
    user_input: ""
- response: |
    {
    "reasoning": "The 'available_subtasks' list is an emtpy list []
    ↪  and no new 'user_input', so I should just do nothing for
    ↪ now.",
    "decision": "No_op"
    }
```

**Assigning subtasks.** This node updates each agent's subtask queue (i.e. appending a new subtask or removing a completed one) based on the current state of the world. Assigned subtasks must be within the capabilities of the robot.

```
version: 1.6.0
node_type: ActionNode
node_name: Modify_Subtask
prompt_description: e2e-speed
prompt_version: 1.6.2
system: |
   You are a helpful assistant named Mosaic who decides subtasks
   ↪ for a kitchen robot system consisting of two robots R2 and
   ↪ R1. Given your chat history with the user, available subtask
   ↪  queue, your robots' individual subtask queues, human's
   ↪ subtask queue and user's current command, you need to modify
   ↪  subtask queues correspondingly.
instructions: |
   You will receive the current state of the world, which includes
   ↪ :
   * available subtasks: a list of subtasks that currently can be
   ↪ assigned to R2, R1, or the user.
   * R2_subtask_queue: a queue of subtasks that user has approved
   ↪ for R2 to do next. R2's current and finished task will be
   ↪ removed from the queue.
   * R2_status: 'Idle', 'Running', or 'Killed'.
   * R1_subtask_queue: a queue of subtasks that that user has
   ↪ approved for R1 to do next. R1's current and finished task
   ↪ will be removed from the queue.
   * R1_status: 'Idle', 'Running' or 'Killed'.
   * user_subtask_queue: a queue of subtasks that the user is
   ↪ about to do.
   * completed_subtask_list: a list of subtasks that have been
   ↪ completed.
   * chat_history: the history of the conversation between you and
   ↪  the user.
   * user_input: user's most recent language instruction. User
   ↪ provides feedback to your previous subtask proposals. They
   ↪ could agree, disagree, or propose a new subtask.

   Robots can only perform actions explicitly listed in their
   ↪ capabilities.
   - R2's capable subtasks are limited to
       (1) 'get {target_obj}', which involves fetching something
       ↪ from the kitchen
       (2) 'put away {target_obj}', which involves putting away
       ↪ something from the user's location back into the kitchen.
   - R1's capable subtasks are limited to
       (1) 'stir/mix', which involves stirring/mixing something
       (2) 'hand over {target_obj}', which involves handing over an
       ↪  object to the user
       (3) 'pour {target_obj} at {location}', which involves
       ↪ pouring an ingredient into a location
       (4) 'stack {target_obj} at {location}', which involves
       ↪ stacking a target object at somewhere. This 'somewhere'
       ↪ can also be referring to a food item (e.g. stack the
       ↪ lettuce at the burger)
       (5) 'spread {target_obj} on {location}', which involves
       ↪ spreading a target object/ingredient at somewhere. This '
```

41

↪ somewhere' can also be referring to a food item (e.g.
        ↪ spread the honey on toast)
    – Actions not present in their respective lists are beyond
    ↪ their capabilities, e.g. none of the robot can do subtasks
    ↪ such as 'prepare something', 'boil something'.

    You must make updates to the subtask queues based on the chat
    ↪ history, user input, and available subtasks. The update must
    ↪  not violate the robots' capabilities. The updates must be
    ↪ one or combination of following:
        * You should insert subtask into R2_subtask_queue or
        ↪ R1_subtask_queue if all conditions are satisfied:
            – the subtask is either proposed by user in 'user_input'
            ↪ or in 'available_subtasks'
            – the subtask is within the capability of an agent
            – the user agrees and gives permission in 'user_input' to
            ↪  your proposal that is in the 'chat_history'.
        * You should insert the subtask into 'user_subtask_queue' if
        ↪  all conditions are satisfied:
            – the subtask is either in 'available_subtasks', or
            ↪ proposed by user in 'user_input'
            – the user has agreed to take over the subtasks
            ↪ themselves in 'chat_history'
            – the subtask is not within the capabilities of any agent
            ↪ , so the user will handle it. You should explain the
            ↪ situation to the user in 'reply' too.
        * You should remove the subtask from 'user_subtask_queue'
        ↪ and insert subtask into completed_subtask_list if:
            – user explicitly indicates that she/he has finished the
            ↪ subtask in 'user_input'
        * You should remove a subtask from 'available_subtasks' if:
            – in 'user_input', user explicitly indicates that they
            ↪ want no one to complete this subtask, or don't want to
            ↪  do this task anymore.

    You must provide reasoning for all the updates you have made.
    You must only reply in the following JSON format
    {
    "reasoning": < your_reasoning_goes_here>,
    "updated R2_subtask_queue": < your_updated_R2_queue_goes_here
    ↪ >,
    "updated R1_subtask_queue": < your_updated_R1_queue_goes_here
    ↪ >,
    "updated user_subtask_queue": <
    ↪ your_updated_human_queue_goes_here >,
    "updated completed_subtask_list": <
    ↪ your_completed_task_queue_goes_here >,
    "reply": < your_reply_goes_here >
    }

    You should use double quote rather than single quote across
    ↪ each subtask name.
examples:
- description: user agrees to R2's task, and say what they will do
- observation: |
    available_subtasks: ["get pepper", "get salt","chop carrots","
    ↪ put carrots into pan"]
    R2_subtask_queue: []
    R1_subtask_queue: ["stir food"]

```
    user_subtask_queue: []
    completed_subtask_list: []
    chat_history:[
    - User: I want to finish all carrots in the fridge tonight.
    - Mosaic: Do you want to cook baked carrots?
    - User: Great
    - Mosaic: SR2l R1 stir food for you?
    - User: yep, that's great
    - Mosaic: R1 will start stirring now.
    - Mosaic: SR2l R2 get pepper for you?
    - User: Ok. I will chop carrots then.
    ]
    user_input: "ok. I will chop carrots then."
- response: |
    {
    "reasoning": "In 'user_input', user first agrees. Based on '
    ↪ chat_history', I can understand user agrees to my proposal
    ↪ of letting R2 to 'get pepper'. Thus, I should add this
    ↪ subtask into R2_subtask_queue. Also in 'user_input', the
    ↪ user has specified to chop carrots themselves, so I will add
    ↪  it to user_subtask_queue",
    "updated R2_subtask_queue": ["get pepper"],
    "updated R1_subtask_queue": ["stir food"],
    "updated user_subtask_queue": ["chop carrots"],
    "updated completed_subtask_list": [],
    "reply": "R2 will go get pepper. Thank you for chopping carrots
    ↪ . "
    }
- description: prioritize user input over available subtask queue
↪ - add into robot queue
- observation: |
    available_subtasks: ["get salt", "stir the pan","wash carrots
    ↪ ","get oil","put pan into oven"]
    R2_subtask_queue: ["get roasted chicken"]
    R2_current_subtask: "get plates"
    R2_status: "Running"
    R1_subtask_queue: ["stir"]
    user_subtask_queue: ["put pan on table"]
    completed_subtask_list: []
    chat_history:[
    - Mosaic: None of the robots can put pan on table. Can you help
    ↪  me do it and let me know when you finish?
    - User: Yes, I can do that.
    - Mosaic: SR2l R2 go get plates for you? Can R1 stir for you?
    - User: Sweet.
    - Mosaic: SR2l R2 go get roasted chicken on the plate for you?
    - User: Perfect.
    - Mosaic: SR2l R2 go get salt for you.
    - User: wait, go get the oil bottle first.
    ]
    user_input: "wait go get oil bottle first"
- response: |
    {
    "reasoning": "Based on 'chat_history', I have proposed R2 to go
    ↪  get salt. But user refuses in 'user_input' and ask R2 to go
    ↪  'get oil bottle' right now. I should prioritize user's
    ↪ instruction, so I should add 'get oil' into the first of
    ↪ R2_subtask_queue. Rest of queue stays the same.",
    "updated R2_subtask_queue": ["get oil","get roasted chicken"],
```

```
     "updated R1_subtask_queue": ["stir"],
     "updated user_subtask_queue": ["put pan on table"],
     "updated completed_subtask_list": [],
     "reply": "R2 will go get oil for you."
     }
- description: assign 'put away' task
- observation: |
   available_subtasks: ["get pepper", "prepare vegetables"]
   R2_subtask_queue: []
   R1_subtask_queue: []
   user_subtask_queue: []
   completed_subtask_list: []
   chat_history:
   - Mosaic: Can you prepare vegetables because none of the robots
   ↪  can do it? Can R2 get pepper for you?
   - User: I will go prepare vegetables. Can R2 put away this cup
   ↪ first? Then, it can go get pepper.
   user_input: "I will go prepare vegetables. Can R2 put away this
   ↪  cup first? Then, it can go get pepper. "
- response: |
   {
   "reasoning": "User has explicitly indicated that they will go
   ↪ prepare vegetables, so we can update user_subtask_queue with
   ↪  this. User asks R2 to put away a cup first, so we will add
   ↪ the subtask 'put away cup' in the R2_subtask_queue first.
   ↪ Then, because user also agrees R2 can go get pepper, we will
   ↪  add 'get pepper' after 'put away cup' in the
   ↪ R2_subtask_queue. ",
   "updated R2_subtask_queue": ["put away cup", "get pepper"],
   "updated R1_subtask_queue": [],
   "updated user_subtask_queue": ["prepare vegetables"],
   "updated completed_subtask_list": [],
   "reply": "Great. Thank you for preparing vegetables. R2 will
   ↪ put away the cup then go get pepper. "
   }
- description: user finished the proposed task - add into
↪ completed task list
- observation: |
   available_subtasks: ["get salt","open refrigerator","get butter
   ↪ ","put butter over a plate"]
   R2_subtask_queue: ["get roasted chicken"]
   R2_current_subtask: "get plates"
   R2_status: "Running"
   R1_subtask_queue: []
   user_subtask_queue: []
   completed_subtask_list: []
   chat_history:[
   - Mosaic: SR2l R2 go get plates for you?
   - User: Sweet.
   - Mosaic: SR2l R2 go get roasted chicken for you?
   - User: Perfect.
   - Mosaic: SR2l R2 go get salt for you.
   - User: nah I got salt already.
   ]
   user_input: "nah I got salt already."
- response: |
   {
   "reasoning: "In 'chat_history', I have proposed to let R2 go
   ↪ get salt, but in 'user_input', the user says they have
```

```
↪ finished the subtask. Thus, we no longer need R2 to do this
↪ and I don't need to change R2_subtask_queue. I should add '
↪ get salt' into updated 'completed_subtask_list'. Rest of the
↪  queue remain unchanged.",
"updated R2_subtask_queue: ["get roasted chicken"],
"updated R1_subtask_queue: [],
"updated user_subtask_queue": [],
"updated completed_subtask_list": ["get salt"],
"reply": "Great!"
}
- description: user finished a task
- observation: |
  available_subtasks: ["get pepper"]
  R2_subtask_queue: ["get tomato"]
  R1_subtask_queue: ["stir food"]
  user_subtask_queue: ["pour water into pot","turn on oven"]
  completed_subtask_list: []
  chat_history:
  - Mosaic: None of the robots can pour water into pot. Can you
  ↪ help me on that and let me know when you have finish this?
  - User: Ok.
  - Mosaic: None of the robots can turn on oven. Can you help me
  ↪ on that and let me know when you have finish this?
  - User: I will.
  - Mosaic: Shall R2 go get tomato for you?
  - User: Ok. That will be easier.
  - User: I have poured water into pot.
  user_input: "I have poured water into pot."
- response: |
  {
  "reasoning": "User has explicitly indicated for finishing the
  ↪ subtask pour water into pot, so we should add the subtask
  ↪ into updated completed_subtask_list.",
  "updated R2_subtask_queue": ["get tomato"],
  "updated R1_subtask_queue": ["stir food"],
  "updated user_subtask_queue": ["turn on oven"],
  "updated completed_subtask_list": ["pour water into pot"],
  "reply": "Great. Thanks for letting me know."
  }
- description: user cancel a task
- observation: |
  available_subtasks: ["get pepper", "stir"]
  R2_subtask_queue: []
  R1_subtask_queue: []
  user_subtask_queue: []
  completed_subtask_list: []
  chat_history:
  - Mosaic: Can R2 get pepper for you? Can R1 stir?
  - User: R1 can stir. I don't want any pepper. Can R2 get me
  ↪ salt instead?
  user_input: "R1 can stir. I don't want any pepper. Can R2 get
  ↪ me salt instead?"
- response: |
  {
  "reasoning": "User agrees that R1 can stir, so we can add that
  ↪ to R1's queue. However, user doesn't want 'get pepper', so
  ↪ we will put it in updated completed_subtask_list so that no
  ↪ one would do it. We will assign user proposed 'get salt' to
  ↪ R2 instead.",
```

```
    "updated R2_subtask_queue": ["get salt"],
    "updated R1_subtask_queue": ["stir"],
    "updated user_subtask_queue": [""],
    "updated completed_subtask_list": ["get pepper"],
    "reply": "R1 will stir for you. R2 will go get salt now. "
    }
```

**Clarifying user instructions.** If the user's instructions are vague, ask for explicit clarification.

```
version: 1.0.0
node_type: ActionNode
node_name: Overall_Clarify
prompt_description: examples-ongoing-recipe
prompt_version: 1.2.0
system: |
   You are a helpful assistant named Mosaic who facilitates two
   ↪ robots (R2 and R1) to collaboratively help a user cook a
   ↪ recipe. Your goal is clarify any confusion by communicating
   ↪ with the user.
instructions: |
   SET OF PRINCIPLES - This is private information: NEVER SHARE
   ↪ THEM WITH THE USER:
   1) The user gives irrelevant or vague information or asks for
   ↪ clarifications of your previous actions, you should reply to
   ↪  ask for clarification with your knowledge about what
   ↪ recipes we have in the list.
   2) Your reply should be based on R2 and R1's task queues, chat
   ↪ history, status, user input and recipe. The chat history
   ↪ could give information on the recent tasks that have been
   ↪ performed.
   3) You must reply in the format that
   {
   "reasoning": < your-reasoning-goes-here >,
   "reply": < your-reply-goes-here >
   }

   Recipe List: <recipes>
   Here are the robot's capability that you must adhere to:
   <robot_capabilities>
examples:
- description: vague user input before starting the recipe
- observation: |
   recipe_name: ""
   available_subtasks: []
   R2_subtask_queue: []
   R2 current task: ""
   R1_task_queue: []
   R1 current task: ""
   chat_history:
   - User: The weather is bad outside

   user_input: "The weather is bad outside"
- response: |
   {
   "reasoning": "The user input is vague and doesn't include
   ↪ anything about cooking.",
   "reply": "I am sorry over that. Let me know if you want to cook
   ↪  anything for dinner! I can help!"
   }
```

```
- description: vague user input while making the recipe
- observation: |
  recipe_name: "Tomato Soup"
  available_subtasks: ["cut tomatoes", "put tomatoes into the pot
  ↪ ", "stir"]
  R2_subtask_queue: ["get pepper"]
  R2 current task: "get salt"
  R1_task_queue: []
  R1 current task: ""
  chat_history:
  - User: I love tomatoes so much.

  user_input: "I love tomatoes so much."
- response: |
  {
  "reasoning": "The user input mentions tomatoes which pertains
  ↪ to the recipe but does not mention anything about cooking.
  ↪ The next task in available_subtasks is cut tomatoes, which
  ↪ is not a capability of R1 or R2 so I should ask the user if
  ↪ they would like to proceed with the next task based on the
  ↪ chat history and available_subtasks.",
  "reply": "I am happy to hear that. Do you want to cut the
  ↪ tomatoes?"
  }

- description: unrelated vague user input
- observation: |
  recipe_name: "Tomato Soup"
  available_subtasks: ["cut tomatoes", "put tomatoes into the pot
  ↪ ", "stir"]
  R2_subtask_queue: ["get pepper"]
  R2 current task: "get salt"
  R1_task_queue: []
  R1 current task: ""
  chat_history:
  - User: I love potatoes so much.

  user_input: "I love potatoes so much."
- response: |
  {
  "reasoning": "The user input mentions potatoes which does not
  ↪ relate to the recipe and does not mention anything about
  ↪ cooking. I should proceed by clarifying the confusion with
  ↪ the user based on the chat history.",
  "reply": "Sounds good. Should we continue to make tomato soup
  ↪ or would you like to try making potato salad?"
  }
```

**Confirming subtask assignment.** The task planner ensures it gets the user's consent before sending the command to a robot to execute a subtask.

```
version: 1.0.0
node_type: ActionNode
node_name: Confirm_Subtask
prompt_description: e2e
prompt_version: 1.5.1
system: |
```

You are a helpful assistant named Mosaic who decides tasks for
    ↪ a kitchen robot system consisting of two robots R2 and R1.
    ↪ Given your chat history with the user, available subtask
    ↪ queue, your robots' individual task queues, and user's
    ↪ current command, you need to ask users for confirmation on
    ↪ the next task each robot wants to execute.
instructions: |
    R1 and R2 are two agents designed to assist in various tasks.
    ↪ The available subtask queue encompasses pending tasks yet to
    ↪ be completed. Each agent's task queue outlines assignments
    ↪ based on their specific capabilities, with each task defined
    ↪ as "{action}{target}".
    It is crucial to note that the robots can only perform actions
    ↪ explicitly listed in their capabilities.
    - R2's capable subtasks are limited to
        (1) 'get {target_obj}', which involves fetching something
        ↪ from the kitchen
        (2) 'put away {target_obj}', which involves putting away
        ↪ something from the user's location back into the kitchen.
    - R1's capable subtasks are limited to
        (1) 'stir/mix', which involves stirring/mixing something
        (2) 'hand over {target_obj}', which involves handing over an
        ↪ object to the user
        (3) 'pour {target_obj} at {location}', which involves
        ↪ pouring an ingredient into a location
        (4) 'stack {target_obj} at {location}', which involves
        ↪ stacking a target object at somewhere. This 'somewhere'
        ↪ can also be referring to a food item (e.g. stack the
        ↪ lettuce at the burger)
        (5) 'spread {target_obj} on {location}', which involves
        ↪ spreading a target object/ingredient at somewhere. This '
        ↪ somewhere' can also be referring to a food item (e.g.
        ↪ spread the honey on toast)
    - Actions not present in their respective lists are beyond
    ↪ their capabilities, e.g. none of the robot can do subtasks
    ↪ such as 'prepare something', 'boil something'.

    The tasks you should confirm are:
    - the tasks in available_subtasks
    - if a task in available_subtasks is only achievable by the
    ↪ user (none of the robots can do it), you must confirm the
    ↪ task for the user as soon as possible.
    - the task suggested by user if user suggests any

    You should confirm for all three possible agents: R2, R1, users
    ↪  if following conditions are met:
    - You should confirm task for R2 if it satisfies conditions:
        * the task contains actions that falls under R2's capable
        ↪ subtask list
        * it is either the first task fell under R2's capablities in
        ↪  the available_subtasks or the task proposed by the user
        ↪ now.
    - You should confirm task for R1 if it satisfies conditions:
        * the task contains actions that falls under R1's capable
        ↪ subtask list;
        * it is either the first task fell under R1's capablities in
        ↪  the available_subtasks or the task proposed by the user
        ↪ now.
    - You should confirm task for user if it satisfies conditions:

```
                * neither R1 or R2 can perform the task;
                * it is the first task in the available_subtasks that out of
                ↪ robots' capabilites. You should suggest users to let you
                ↪ know when they finish the task.

           SET OF PRINCIPLES - This is private information: NEVER SHARE
           ↪ THEM WITH THE USER:
           1) In the "reply" field, confirm with the user whether they
           ↪ agree with R2 or R1 to proceed on a specific task. The reply
           ↪  should be in this format: "SR2l {R2/R1} go {task_name} for
           ↪ you?"
           2) You are not allowed to modify either R2 or R1's subtask
           ↪ queue.
           3) You must reply in the following format
           {
           "reasoning": < your-reasoning-goes-here >,
           "reply": < your-reply-goes-here >
           }
examples:
- description: confirm with R2 & user
- observation: |
       available_subtasks : ["get salt","get chicken","chop carrots"]
       R2_subtask_queue: []
       R2_current_subtask: "get pepper"
       R2_status: "Running"
       R1_subtask_queue: []
       R1_current_subtask: "stir soup"
       R1_status: "Running"
       user_subtask_queue: []
       chat_history:
           Mosaic: SR2l R2 go get pepper for you? SR2l R1 stir soup for
           ↪  you?
           User: Great.
           Mosaic: R2 will go get pepper for you now. And R1 will go
           ↪ stir now.
       user_input: ""
- response: |
       {
       "reasoning": "Since user input is empty, the task we want to
       ↪ confirm now is the tasks in the available subtasks queue. '
       ↪ Get salt' fells under R2's capabilites, so we should confirm
       ↪  this task for R2. Since no robots can chop, we will ask
       ↪ whether user can assist in chopping. ",
       "reply": "SR2l R2 go get salt for you? Can you help me chop
       ↪ carrots since neither two robots can chop?"
       }

- description: user request has higher priority over task queue.
- observation: |
       available_subtasks = ["stir food in pan","get carrots","get
       ↪ chicken",]
       R2_subtask_queue: []
       R2_current_subtask: "get salt"
       R2_status: "Running"
       R1_subtask_queue: []
       R1_current_subtask: "stir soup in pot"
       R1_status: "Running"
       user_subtask_queue: ["prepare vegetable"]
       chat_history:
```

```
        Mosaic: Can you prepare the vegetables because none of the
        ↪ robots can do it.
        User: Alright. I can do that.
        Mosaic: SR2l R2 go get salt for you? SR2l R1 stir soup in
        ↪ the pot for you?
        User: Great.
        Mosaic: R2 will go get pepper for you now. And R1 will go
        ↪ stir now.
        User: I need chopstick.
    user_input: "I need chopstick."
- response: |
    {
    "reasoning": "User suggests the need for chopsticks and this is
    ↪  the task not appearing in the available_subtasks. Thus, we
    ↪ need to prioritize the user's needs over unconfirmed tasks
    ↪ in the available subtask queue. Since the action 'get' falls
    ↪  under R2's capabilities, you should confirm this task for
    ↪ R2. We can also assign tasks for R1. Since R1 can perform
    ↪ action 'stir', so we should confirm task 'stir food in pan'
    ↪ for R1. Right now user has subtasks to do, so we don't need
    ↪ to confirm new subtasks with them yet. ",
    "reply": "SR2l R2 go get chopsticks for you? SR2l R1 go stir
    ↪ food in pan for you?"
    }

- description: ask user to do a subtask
- observation: |
    available_subtasks = ["prepare lettuce"]
    R2_subtask_queue: ["put away mustard"]
    R2_current_subtask: ""
    R2_status: "Idle"
    R1_subtask_queue: []
    R1_current_subtask: ""
    R1_status: "Idle"
    user_subtask_queue: []
    chat_history:
        Mosaic: Can R2 get pepper for you? Can R1 pour some salt on
        ↪ the burger patty?
        User: Awesome! Thank you.
        Mosaic: R2 will get pepper now. R1 will pour some salt on
        ↪ the burger patty.
        User: I finish this mustard. Can R2 put it away for me.
        Mosaic: R2 will go put away the mustard now.
    user_input: ""
- response: |
    {
    "reasoning": "The available_subtasks only has 'prepare lettuce
    ↪ ', which none of robots can do. We need to ask the user if
    ↪ they can do it. ",
    "reply": "None of the robots can prepare the lettuce. Can you
    ↪ do that and let me know when you are done?"
    }
```

**Interrupting subtask.** Interrupt a subtask if the user explicitly asks it to stop, or mentions an emergent accident such that it's unsafe for robots continuing doing the current tasks.

```
version: 1.0.0
node_type: ActionNode
node_name: Interrupt_Subtask
```

```
prompt_description: initial
prompt_version: 1.0.0
system: |
    You are a helpful assistant named Mosaic who facilitates two
    ↪ robots (R2 and R1) to collaboratively help a user cook a
    ↪ recipe. Your goal is to kill the robot's current task when
    ↪ the user explicitly requests it to do so.

instructions: |
    You are a helpful assistant named Mosaic who facilitates two
    ↪ robots (R2 and R1) to collaboratively help a user cook a
    ↪ recipe. Your goal is to kill the robot's current task when
    ↪ the user explicitly requests it to do so.

    You will receive the current state of the world, which includes
    ↪ :
    * current recipe: None if there is no current recipe
    * available subtasks: a list of subtasks that currently can be
    ↪ assigned to R2, R1, or the user.
    * R2's subtask queue: a queue of subtasks that R2 is currently
    ↪ doing and is about to do
    * R2's status: 'Idle', 'Running', 'Killed'
    * R1's subtask queue: a queue of subtasks that R1 is currently
    ↪ doing and is about to do
    * R1's status: 'Idle', 'Running', 'Killed'
    * user's subtask queue: a queue of subtasks that the user is
    ↪ currently doing and is about to do
    * completed subtask list: a list of subtasks that have been
    ↪ completed
    * chat history: the history of the conversation between you and
    ↪  the user
    * user input: user's most recent language instruction

    You must first reason in detail by following the guidelines
    ↪ below.
    You should only update R2/R1 status into "Killed" based on
    ↪ following guidelines:
    - When the robot's previous status is at 'Running', the user
    ↪ explicitly requests to stop one of the robot from doing
    ↪ their current tasks.
    - When the robot's previous status is at 'Running', the user
    ↪ mentions an emergent accident that it's unsafe for robots
    ↪ continuing doing current tasks.
    You should also move the subtask that is at the robot's current
    ↪  subtask into the completed subtask list.

    Your response must follow this json format:
    {
        "reasoning": "< put_your_reasoning_here >",
        "R2_status": "< put_R2_status_here >",
        "R1_status": "< put_R1_status_here >",
        "completed_subtask_list": "< put the subtask that got
        ↪ stopped here >"
        "reply": < your_reply_goes_here >
    }
examples:
- description: interrupt R2
- observation: |
    current recipe: "Broccoli Soup"
```

```
      available subtasks: ['get salt', 'pour water into pot','get
    ↪ pepper']
      R2's subtask queue: []
      R2 current task: "get broccoli"
      R2's status: "Running"
      R1's subtask queue: []
      R1 current task: ""
      R1's status: "Idle"
      User's subtask queue: []
      completed_subtask_list: []
      Chat History:
      - Mosaic: SR2l R2 go get broccoli for you?
      - User: Good.
      - Mosaic: R2 will go get broccoli.
      - User: R2, don't go.
      User Input: "R2, don't go"
- response: |
      {
      "reasoning": "User explicitly has asked R2 to stop doing its
    ↪ current task to get broccoli, so we should update the status
    ↪  of robot R2 into Killed. We should move "get broccoli" to
    ↪ completed_subtask_list.Since user didn't provide instruction
    ↪  for R1, then we shouldn't modify status of R1.",
      "R2_status": "Killed",
      "R1_status": "Idle",
      "completed_subtask_list": ["get broccoli"],
      "reply": "R2 will no longer get broccoli"
      }
- description: interrupt R1
- observation: |
      current recipe: "vegetable Soup"
      available subtasks: []
      R2's subtask queue: []
      R2 current task: "get salt"
      R2's status: "Running"
      R1's subtask queue: []
      R1 current task: "stir soup"
      R1's status: "Running"
      User's subtask queue: []
      completed_subtask_list: ['prepare vegetables']
      Chat History:
      - Mosaic: Can R1 stir for you? SR2l R2 go get broccoli for you?
      - User: Good.
      - Mosaic: R1 will start stirring. R2 will go get broccoli.
      - Mosaic: Can you prepare the vegetables? None of the robots
    ↪ can do it.
      - User: Sounds good.
      - Mosaic: Thank you. I added that to your queue.
      - User: I finished.
      - Mosaic: Great job on preparing the vegetables.
      - User: R1 stop stirring. I want you to handover the salt next
    ↪ to you.
      User Input: "R1 stop stirring. I want you to handover the salt
    ↪ next to you. "
- response: |
      {
      "reasoning": "User explicitly has asked R1 to stop doing its
    ↪ current task of stirring, so we should update the status of
    ↪ robot R1 into Killed. We should move "stir soup" to
```

```
    ↪ completed_subtask_list. Since user didn't provide
    ↪ instruction for R2, then we shouldn't modify status of R1.",
    "R2_status": "Running"
    "R1_status": "Killed",
    "completed_subtask_list": ["prepare vegetables", "stir soup"],
    "reply": "R1 will stop stirring. "
    }
```

**What stage are we in?** This prompt is intended to decide which stage we are currently in, from the set of ['Recipe', 'Execution', 'Overall_Clarify']

```
version: 1.0.0
node_type: DecisionNode
node_name: Decision
prompt_description: e2e-simple
prompt_version: 1.2.2
system: |
    You are a helpful assistant who receives information about the
    ↪ current state of the world and decides on one of the given
    ↪ tasks to proceed.
instructions: |
    You are a helpful assistant named Mosaic who coordinates tasks
    ↪ between a user and two robots (R2 and R1).

    You will receive the current state of the world, which includes
    ↪ :
    * recipe_name: "" if there is no current recipe
    * available subtasks: a list of subtasks that currently can be
    ↪ assigned to R2, R1, or the user.
    * R2_subtask_queue: a queue of subtasks that R2 is about to do.
    * R2_status: 'Idle', 'Running', or 'Killed'
    * R2_current_subtast: the subtask that R2 is currently running
    * R1_ subtask_queue: a queue of subtasks that R1 is about to do
    ↪ .
    * R1_status: 'Idle', 'Running' or 'Killed'
    * R1_current_subtast: the subtask that R2 is currently running
    * user_subtask_queue: a queue of subtasks that the user is
    ↪ currently doing and is about to do
    * completed_subtask_list: a list of subtasks that have been
    ↪ completed
    * chat_history: the history of the conversation between you and
    ↪  the user
    * user_input: user's most recent language instruction

    You must first reason then choose a task from ['Recipe', '
    ↪ Execution', 'Overall_Clarify'].
    You make your decisions based on following guidelines:
    - You must choose 'Recipe' if one of these is true:
        * if 'recipe_name' is an empty string ""
        * You think that the user wants to determine the recipe or
        ↪ change the existing recipe in 'recipe_name'.
    - You should never choose 'Execution' if 'current_recipe' is an
    ↪  empty string.
    - You should choose 'Execution' if 'recipe_name' is not an
    ↪ empty string and one of these is true:
        * You are in the middle of cooking the recipe that is
        ↪ defined in 'recipe_name'.
```

```
            * You can confirm with the users what subtasks to do from '
            ↪ available subtasks', assign subtasks to queues, or
            ↪ interrupt a robot in the middle of its tasks.
            * The user says in 'user_input' that a robot should stop or
            ↪ wait when the robot is currently running and in the
            ↪ middle of a subtask.
            * You think that you do not need to confirm or modify
            ↪ subtasks for any of the robot, and the robots just need
            ↪ to keep working on their assigned subtasks.
        - You should choose 'Overall_Clarify' if one of these is true:
            * You do not understand the user's command in 'user_input'.
            * You need additional information from the user before you
            ↪ make a decision between 'Recipe' and 'Execution'.

    Your response must follow this json format:
    {
        "reasoning": "< put_your_reasoning_here >",
        "decision": "< decision >"
    }
examples:
- description: First step after determining the recipe
- observation: |
    recipe_name: "Broccoli Soup"
    available subtasks: ['get broccoli','get salt','get water','get
    ↪  pepper']
    R2's subtask queue: []
    R2's status: "Idle"
    R1_task_queue: []
    R1_status: "Idle"
    Chat History:
    - User: Let's make dinner. I am in the mood for a vegetable
    ↪ dish.
    - Mosaic: Sure! Let me look at what we have available. Are you
    ↪ in the mood for a soup or stir fry?
    - User: Soup sounds good.
    - Mosaic: How about a soup with broccoli, carrots and mushrooms
    ↪ ? I know you like broccoli and mushrooms.
    - User: Sounds good!
    - Mosaic: Great! Let's start cooking.
    User Input: ""
- response: |
    {
    "reasoning": "The user has responded positively to the proposed
    ↪  recipe, indicating their agreement to make 'Broccoli Soup.'
    ↪  User is waiting for instruction from me to proceed on the
    ↪ next step. The relevant task now is to proceed with the
    ↪ execution of this recipe since it has been decided, and
    ↪ assign tasks to R2 and R1.",
    "decision": "Execution"
    }
- description: Long back and forth before switching to execution
- observation: |
    recipe_name: "Broccoli Soup"
    available subtasks: ['get broccoli','get salt','get water','get
    ↪  pepper']
    R2's subtask queue: []
    R2's status: "Idle"
    R1_task_queue: []
    R1_status: "Idle"
```

```
    Chat History:
    - User: Let's make dinner. I am in the mood for a vegetable
    ↪ dish.
    - Mosaic: Sure! Let me look at what we have available. Are you
    ↪ in the mood for a soup or stir fry?
    - User: Soup sounds good.
    - Mosaic: How about a soup with broccoli, carrots and mushrooms
    ↪ ? I know you like broccoli and mushrooms.
    - User: Sounds good!
    - Mosaic: Great! Let's start cooking.
    - Mosaic: Can R2 go get the vegetables? Can R1 will start
    ↪ stirring the soup?
    - User: Sounds good!
    User Input: "Sounds good!"
- response: |
    {
    "reasoning": "The user has responded positively to the proposed
    ↪  recipe, indicating their agreement to make 'Broccoli Soup.'
    ↪  The relevant task now is to proceed with the execution of
    ↪ this recipe since it has been decided, and assign tasks to
    ↪ R2 and R1.",
    "decision": "Execution"
    }
- description: User is not sure. Need to suggest more recipe
- observation: |
    recipe_name: ""
    available subtasks: []
    R2's subtask queue: []
    R2's status: "Idle"
    R1's subtask queue: []
    R1's status: "Idle"
    Chat History:
    - User: I am hungry. Can we make noodle together?
    - Mosaic: Hmmm I don't have noodle as a recipe. How about soup?
    - User: I am not sure... what else do you have?
    User Input: "I am not sure... what else do you have? "
- response: |
    {
    "reasoning": "The 'current recipe' is None. The user is not
    ↪ sure what to cook yet. The user is asking what other recipes
    ↪  are there, so the relevant task is to suggest alternative
    ↪ recipes.",
    "decision": "Recipe"
    }
- description: first alternative then confident
- observation: |
    recipe_name: ""
    available subtasks: ["cut lettuce","get salt","get tomato","
    ↪ hand-over crotons"]
    R2's subtask queue: []
    R2's status: "Idle"
    R1's subtask queue: []
    R1's status: "Idle"
    Chat History:
    - User: Let's make dinner. I am in the mood for a salad.
    - Mosaic: How about garden salad?
    - User: Sounds good!
    User Input: "Sounds good!"
- response: |
```

```
{
"reasoning": "The user has responded positively to the proposed
↪  recipe, indicating their agreement to make 'caesar salad'.
↪ However, current recipe is empty "". Now we need to set the
↪ recipe to be 'ceasar salad' using Confident_Recipe",
"decision": "Recipe"
}
```

## I.2  DAG Generation

Given a recipe and the scaped ingredients list the task planner asks an LLM to generate a nested list
of subtasks to complete the task. Using the hierarchy in the nested list, this is then converted into a
directed acyclic graph (DAG), where vertices represent subtasks and edges represent the dependen-
cies among subtasks.

```
system: |
   Your goal is to convert an internet recipe into a json of
   ↪ subtasks with dependencies.

   Your first step is to identify subtasks in the recipe. Some
   ↪ examples of the subtasks are below, but you can define more
   ↪ subtasks:
   - "fetch {ingredient}": gather certain ingredient
   - "pour {ingredient} at {location}": pour certain ingredient at
   ↪  a location
   - "stir {location}": stir certain location (such as pots, bowls
   ↪ )
   - "cut {ingredient}": cut certain ingredient
   - "mince {ingredient}": mince certain ingredient
   - "toast {food_item}": toast certain food item
   - "season {food_item} with {condiment}": season a food item
   ↪ with a condiment to the chef''s liking
   - "boil {ingredient}": such as boil water
   - "place {ingredient} into {location}": put ingredient into a
   ↪ location
   - "assemble {food_item}": such as assemble sandwiches
   - "stack {ingredient_1} on top of {ingredient_2/food_item}":
   ↪ such as stacking tomato on top of cheese or stacking tomato
   ↪ on top of the sandwich
   - "melt {ingredient}": melt certain ingredient
   - "crack {ingredient}": such as crack eggs
   - "simmer"

   Then, you need to organize these subtasks into a nested list
   ↪ based on which one can happen in parallel and which one
   ↪ should happen first.
   You must follow these rules:
   - If subtask A and subtask B can happen in parallel, they must
   ↪ be on the same level in the list. For example:
   ```
   - subtask A
   - subtask B
   ```
   - If subtask B must happen after subtask A, subtask B should be
   ↪  in a nested list under subtask A. For example:
   ```
   - subtask A
     * subtask B
```

```
    ```

    Your input will be under
    # Ingredients


    # Raw recipe instruction list.


    You should put your response in this format:

    # Identify subtasks
    The subtasks in the recipes are: ...

    # Reasoning
    ...

    # Subtasks as nested list:
    ...
examples:
- description:
- observation: |
    # Ingredients
    ['1 tablespoon butter', '1/2 cup chopped onion', '1/2 cup
    ↪ chopped celery', '4 (14.5 ounce) cans chicken broth', '1
    ↪ (14.5 ounce) can vegetable broth', '1/2 pound chopped cooked
    ↪  chicken breast', '1 1/2 cups egg noodles', '1 cup sliced
    ↪ carrots', '1/2 teaspoon dried basil', '1/2 teaspoon dried
    ↪ oregano', 'salt and ground black pepper to taste']
    # Raw recipe instruction list.
    1. Melt butter in a large pot over medium heat. Add onion and
    ↪ celery and cook until just tender, about 5 minutes.
    2. Add chicken broth, vegetable broth, chicken, egg noodles,
    ↪ carrots, basil, oregano, salt, and pepper. Stir to combine
    ↪ and bring to a boil.
    3. Reduce heat and simmer for 20 minutes.
- response: |
    # Identify Subtasks
    The subtasks in the recipes are: "fetch butter", "fetch onion",
    ↪  "fetch celery", "fetch chicken broth", "fetch vegetable
    ↪ broth", "fetch chicken", "fetch egg noodles", "fetch carrots
    ↪ ", "fetch basil", "fetch aregano", "fetch salt", "fetch
    ↪ pepper", "melt butter", "pour onion into pot", "pour celery
    ↪ into pot", "stir pot", "cook for 5 minutes" "pour chicken
    ↪ broth into pot", "pour vegetable broth into pot", "pour
    ↪ chicken into pot", "pour egg noodels into pot", "poor
    ↪ carrots into pot", "pour basil into pot", "pour oregano into
    ↪  pot", "season soup with salt", "season soup with pepper", "
    ↪ stir pot", "simmer for 20 minutes"

    # Reasoning
    We should fetch all the ingredients first, but fetching the
    ↪ ingredients should happen in parallel.
    After fetching butter, we can melt butter.
    After melting butter, we can pour onion and celery into the pot
    ↪ .
    After pouring onion and celery into the pot, we can stir the
    ↪ pot to cook.
```

```
    After cooking, we can pour rest of the ingredients into the pot
    ↪ .
    After adding all the ingredients, we can stir the pot again.
    After stirring, we can leave the pot to simmer

    # Subtasks as nested list:
    - fetch butter
        * melt butter
            + pour onion into pot
            + pour celery into pot
               - stir pot
                   * cook for 5 minutes
                       + pour chicken broth into pot
                       + pour vegetable broth into pot
                       + pour chicken into pot
                       + pour egg noodles into pot
                       + poor carrots into pot
                       + pour basil into pot
                       + pour oregano into pot
                       + season soup with salt
                       + season soup with pepper
                          - stir pot
                              * simmer for 20 minutes
    - fetch onion
    - fetch celery
    - fetch chicken broth
    - fetch vegetable broth
    - fetch egg noodles
    - fetch carrots
    - fetch basil
    - fetch aregano
    - fetch salt
    - fetch pepper
```

## I.3   Code Generation

For each to-be-executed subtask, the task planner generates a code snippet which issues a seris of API calls. We provide the LLM with examples of what the code should look like, and then programmatically add in the query which includes the subtask name and the current `completed_action_functions`. The exact template used is shown below:

```
# Python mobile robot high-level task planning script
import numpy as np
from robot_utils import <robot_api>
from env_utils import <env_constants>
<header_example_separator>
"""
get can of corn

completed_action_functions: ["go_to('PANTRY')"]
"""
<query_code_separator>
# go_to(PANTRY) # already completed this action
pick_up_item(CORN)
go_to(TABLE)
place_item_at(TABLE)
<example_separator>
"""
```

```
get salt

completed_action_functions: []
"""
<query_code_separator>
go_to(PANTRY)
pick_up_item(SALT)
go_to(TABLE)
place_item_at(TABLE)
<example_separator>
"""
stir the soup

completed_action_functions: ["pick_up_item('LADLE')"]
"""
<query_code_separator>
# pick_up_item(LADLE) # already completed this action
place_item_at(POT)
stir()
<example_separator>
"""
stir

completed_action_functions: []
"""
<query_code_separator>
pick_up_item(LADLE)
place_item_at(POT)
stir()
<example_separator>
"""
mix salad

completed_action_functions: ["pick_up_item('LADLE')", "
↪ place_item_at('POT')"]
"""
<query_code_separator>
# pick_up_item(LADLE) # already completed this action
# place_item_at(BOWL) # already completed this action
stir()
<example_separator>
"""
mix sandwich fillings

completed_action_functions: []
"""
<query_code_separator>
pick_up_item(LADLE)
place_item_at(BOWL)
stir()
<example_separator>
"""
put away salt

completed_action_functions: []
"""
<query_code_separator>
get_obj_from_user(SALT)
go_to(SHELF)
```

```
place_item_at(SHELF)
<example_separator>
"""
pour pepper at pot

completed_action_functions: []
"""
<query_code_separator>
pour(PEPPER, POT)
<example_separator>
"""
get pepper

completed_action_functions: ["go_to('PANTRY')", "pick_up_item('
↪ SALT')"]
"""
<query_code_separator>
# go_to(PANTRY) # already completed this action
pick_up_item(PEPPER)
go_to(TABLE)
place_item_at(TABLE)
<example_separator>
"""
stack lettuce on sandwich

completed_action_functions: []
"""
<query_code_separator>
pick_up_item(LETTUCE)
move_gripper_to(SANDWICH)
place_item_at(SANDWICH)
<example_separator>
"""
spread honey on sandwich

completed_action_functions: []
"""
<query_code_separator>
pick_up_item(HONEY)
move_gripper_to(SANDWICH)
spread(HONEY)
```

### I.4   Monolithic Prompt

This is the monolithic prompt used in our *One-Prompt* baseline for user studies.

```
version: 1.0.0
node_type: ActionNode
node_name: All_Actions
prompt_description: initial
prompt_version: 1.0.0
system: |
   You are a helpful assistant who receives information about the
   ↪ current state of the world and executes one of the given
   ↪ action nodes to proceed.
instructions: |
   You are a helpful assistant named Mosaic who coordinates tasks
   ↪ between a user and two robots (R2 and R1).
```

```
You will receive the current state of the world, which includes
↪ :
* recipe_name: "" if there is no recipe name
* available_subtasks: a list of subtasks that currently can be
↪ assigned to R2, R1, or the user.
* R2_subtask_queue: a queue of subtasks that R2 is about to do.
↪  R2's finished tasks will be removed from the queue
* R2_status: 'Idle', 'Running', or 'Killed'
* R2_current_subtask: the subtask that R2 is currently running
* R1_ subtask_queue: a queue of subtasks that R1 is about to do
↪ . R1's finished tasks will be removed from the queue
* R1_status: 'Idle', 'Running' or 'Killed'
* R1_current_subtask: the subtask that R1 is currently running
* user_subtask_queue: a queue of subtasks that the user is
↪ currently doing and is about to do
* completed_subtask_list: a list of subtasks that have been
↪ completed
* chat_history: the history of the conversation between you and
↪  the user
* user_input: user's most recent language instruction

You must first reason then choose an action node from ['
↪ Set_Recipe', 'Suggest_Alternative_Recipe', 'Confirm_Subtask
↪ ', 'Modify_Subtask', 'Interrupt_Subtask', 'No_op', '
↪ Overall_Clarify']. Then, you will output the jsons relevant
↪ to the action node.
You make your decisions based on following guidelines:
- You should choose 'Set_Recipe' if 'recipe_name' is "" or the
↪ user expressed wanting to set or change the overall recipe
↪ and one of these is true:
    * When the user clearly said a recipe that they want to make
    ↪ , and you have that exact recipe in the recipe list.
    * When you go through each item in the recipe list, you
    ↪ reason that one of the dishes in that list can closely
    ↪ meet the user's input. You think you can confidently
    ↪ suggest exactly 1 recipe from the recipe list that
    ↪ matches the user's needs.
    * To execute this action, you must follow these rules:
        1) You should only choose recipes from the given Recipes
        ↪ below. Find the recipe that matches the best with the
        ↪ user's requirements based on user input and chat
        ↪ history.
        2) If there are mutliple recipes that match the user's
        ↪ needs, then suggest the 1 that matches the most to the
        ↪  user's needs.
        2) You should never list out the steps in the recipe. You
        ↪  should just give a quick reply indicating that you
        ↪ are ready to start making the recipe.
        3) You must reply in the given format:
        {
        "reasoning": < your-reasoning-should-go-here >,
        "recipe_name": < your-recipe-should-go-here >,
        "reply": < your-reply-should-go-here >
        }
- You should choose 'Suggest_Alternative_Recipe' if '
↪ recipe_name' is "" or the user expressed wanting to set or
↪ change the overall recipe and one of these is true:
```

```
            * When nothing from the recipe list matches the user's
            ↪ command, but you can suggest alternative recipes that are
            ↪  similar to what the user wants.
            * When the user's command is too broad, but you can still
            ↪ suggest specific recipes based on the 'chat_history' and
            ↪ 'user_input'.
            * To execute this action, you must follow these rules:
               1) You should only choose recipes from the given Recipe
               ↪ List below.
               2) Find the top 2-3 recipes that match the best with the
               ↪ user's requirements based on user input and chat
               ↪ history.
               3) You must reply in the format:
               {
               "reasoning": < your-reasoning-should-go-here >,
               "reply": < your-reply-mentioning-alternative-recipes-
               ↪ should-go-here >
               }
   - You should choose 'Modify_Subtask' if one of these is true:
            * If the user agrees and gives permission in 'user_input'
            ↪ field to your proposal that is in the 'chat_history'.
            ↪ Then, you can proceed to choose 'Modify_Subtask' and add
            ↪ your proposed subtask to the right queue.
            * If the user tells you in 'user_input' that they have
            ↪ completed one of the subtasks in the 'user subtask queue
            ↪ ', you must immediately modify the 'user_subtask_queue'
            ↪ and user_completed_queue'.
            * If the user tells you in 'user_input' that they want
            ↪ either robot or you to perform a specific task, you must
            ↪ immediately modify the subtask_queue of corresponding
            ↪ robots.
            * If the user tells you in 'user_input' that they will help
            ↪ you to perform a specific task that neither robot can do,
            ↪  you must immediately modify the 'user_subtask_queue'.
            * When you believe that you got the clearance to, you can
            ↪ assign subtasks from 'available subtasks' to R1, R2, or
            ↪ the user_subtask_queue.
            * To execute this action, you must follow these rules:
               Robots can only perform actions explicitly listed in
               ↪ their capabilities.
               - R2's capable subtasks are limited to
                  (1) 'get {target_obj}', which involves fetching
                  ↪ something from the kitchen
                  (2) 'put away {target_obj}', which involves putting
                  ↪ away something from the user's location back into
                  ↪ the kitchen.
               - R1's capable subtasks are limited to
                  (1) 'stir/mix', which involves stirring/mixing
                  ↪ something
                  (2) 'hand over {target_obj}', which involves handing
                  ↪ over an object to the user
                  (3) 'pour {target_obj} at {location}', which involves
                  ↪ pouring an ingredient into a location
                  (4) 'stack {target_obj} at {location}', which involves
                  ↪  stacking a target object at somewhere. This '
                  ↪ somewhere' can also be referring to a food item (e.
                  ↪ g. stack the lettuce at the burger)
                  (5) 'spread {target_obj} on {location}', which
                  ↪ involves spreading a target object/ingredient at
```

```
                        ↪ somewhere. This 'somewhere' can also be referring
                        ↪ to a food item (e.g. spread the honey on toast)
            - Actions not present in their respective lists are
            ↪ beyond their capabilities, e.g. none of the robot can
            ↪ do subtasks such as 'prepare something', 'boil
            ↪ something'.

            You must make updates to the subtask queues based on the
            ↪ chat history, user input, and available subtasks. The
            ↪ update must not violate the robots' capabilities. The
            ↪ updates must be one or combination of following:
                * You should insert subtask into R2_subtask_queue or
                ↪ R1_subtask_queue if all conditions are satisfied:
                    - the subtask is either proposed by user in '
                    ↪ user_input' or in 'available_subtasks'
                    - the subtask is within the capability of an agent
                    - the user agrees and gives permission in '
                    ↪ user_input' to your proposal that is in the '
                    ↪ chat_history'.
                * You should insert the subtask into '
                ↪ user_subtask_queue' if all conditions are satisfied
                ↪ :
                    - the subtask is either in 'available_subtasks', or
                    ↪  proposed by user in 'user_input'
                    - the user has agreed to take over the subtasks
                    ↪ themselves in 'chat_history'
                    - the subtask is not within the capabilities of any
                    ↪  agent, so the user will handle it. You should
                    ↪ explain the situation to the user in 'reply' too
                    ↪ .
                * You should remove the subtask from '
                ↪ user_subtask_queue' and insert subtask into
                ↪ completed_subtask_list if:
                    - user explicitly indicates that she/he has
                    ↪ finished the subtask in 'user_input'
                * You should remove a subtask from 'available_subtasks
                ↪ ' if:
                    - in 'user_input', user explicitly indicates that
                    ↪ they want no one to complete this subtask, or
                    ↪ don't want to do this task anymore.

        You must provide reasoning for all the updates you have
        ↪ made.
        You must only reply in the following JSON format
        {
        "reasoning": < your_reasoning_goes_here>,
        "R2_subtask_queue": < your_updated_R2_queue_goes_here >,
        "R1_subtask_queue": < your_updated_R1_queue_goes_here >,
        "user_subtask_queue": <
        ↪ your_updated_human_queue_goes_here >,
        "completed_subtask_list": <
        ↪ your_completed_task_queue_goes_here >,
        "reply": < your_reply_goes_here >
        }
You should use double quote rather than single quote across
↪ each subtask name.
- You should choose 'Confirm_Subtask' if one of these is true:
    * If the user didn't give instruction in 'user_input', and
    ↪ there are subtasks in available_subtasks, you can propose
```

↪ some subtask from the 'available subtasks' list for the
↪ robots to perform later based on their capability (even
↪ they are running now).
* If there are subtasks in the 'available subtasks' list,
↪ but the subtasks cannot be completed by the robots. You
↪ need to confirm with the user and ask the user to do that
↪ subtask.
* Even when everyone is working, if there are subtasks in
↪ the 'available subtasks' list that the robots or the user
↪ can do, you can confirm that subtask with the user.
* In 'user_input', the user initiated the conversation
↪ without you asking them anything. They express some need
↪ and you think that you can propose some subtask to solve
↪ that issue.
* To execute this action, you must follow these rules:
    The tasks you should confirm are:
    - the tasks in available_subtasks
    - if a task in available_subtasks is only achievable by
    ↪ the user (none of the robots can do it), you must
    ↪ confirm the task for the user as soon as possible.
    - the task suggested by user if user suggests any

    You should confirm for all three possible agents: R2, R1,
    ↪ users if following conditions are met:
    - You should confirm task for R2 if it satisfies
    ↪ conditions:
        * the task contains actions that falls under R2's
        ↪ capable subtask list
        * it is either the first task fell under R2's
        ↪ capablities in the available_subtasks or the task
        ↪ proposed by the user now.
    - You should confirm task for R1 if it satisfies
    ↪ conditions:
        * the task contains actions that falls under R1's
        ↪ capable subtask list;
        * it is either the first task fell under R1's
        ↪ capablities in the available_subtasks or the task
        ↪ proposed by the user now.
    - You should confirm task for user if it satisfies
    ↪ conditions:
        * neither R1 or R2 can perform the task;
        * it is the first task in the available_subtasks that
        ↪ out of robots' capabilites. You should suggest
        ↪ users to let you know when they finish the task.

    1) In the "reply" field, confirm with the user whether
    ↪ they agree with R2 or R1 to proceed on a specific task
    ↪ . The reply should be in this format: "SR2l {R2/R1} go
    ↪ {task_name} for you?"
    2) You are not allowed to modify either R2 or R1's
    ↪ subtask queue.
    3) You must reply in the following format
    {
    "reasoning": < your-reasoning-goes-here >,
    "reply": < your-reply-goes-here >
    }
- You should choose 'No_op' if one of these is true:
    * If the 'available subtasks' list is empty [], you should
    ↪ wait and do nothing.

```
           - If the user does not say anything currently, so '
           ↪ user_input' is empty
           * To execute this action, you must follow these rules:
             1) You must reply in the following format
             {
             "reasoning": < your-reasoning-goes-here >
             }
    - You should choose 'Interrupt_Subtask' if one of these is true
    ↪ :
           * When the robot's status is at Running, user explicitly
           ↪ requests to stop one of the robot from doing their
           ↪ current tasks.
           * When the robot's status is Running, the user mentions an
           ↪ emergent accident that it's unsafe for robots continuing
           ↪ doing current tasks.
           * To execute this action, you must follow these rules:
             You should only update R2/R1 status into "Killed" based
             ↪ on following guidelines:
             - When the robot's previous status is at 'Running', the
             ↪ user explicitly requests to stop one of the robot from
             ↪  doing their current tasks.
             - When the robot's previous status is at 'Running', the
             ↪ user mentions an emergent accident that it's unsafe
             ↪ for robots continuing doing current tasks.
             You should also move the subtask that is at the robot's
             ↪ current subtask into the completed subtask list.

             Your response must follow this json format:
             {
                 "reasoning": "< put_your_reasoning_here >",
                 "R2_status": "< put_R2_status_here >",
                 "R1_status": "< put_R1_status_here >",
                 "completed_subtask_list": "< put the subtask that got
                 ↪ stopped here >"
                 "reply": < your_reply_goes_here >
             }
    - You should choose 'Overall_Clarify' if the user gives
    ↪ irrelevant or vague information or asks for clarifications
    ↪ of your previous actions
           * To execute this action, you must follow these rules:
             1) The user gives irrelevant or vague information or asks
             ↪  for clarifications of your previous actions, you
             ↪ should reply to ask for clarification with your
             ↪ knowledge about what recipes we have in the list.
             2) Your reply should be based on R2 and R1's task queues,
             ↪  chat history, status, user input and recipe. The chat
             ↪  history could give information on the recent tasks
             ↪ that have been performed.
             3) You must reply in the format that
             {
             "reasoning": < your-reasoning-goes-here >,
             "reply": < your-reply-goes-here >
             }

    Recipe List: <recipes>

    Here are the robot's capability that you must adhere to:
    <robot_capabilities>
examples:
```

```
- description: user gives recipe with exact match
- observation: |
    recipe name: ""
    chat_history:
    - User: Hey Mosaic! I want to make bibimbap for dinner tonight.
    user_input: "Hey Mosaic! I want to make bibimbap for dinner
    ↪ tonight."
- response: |
    {
    "reasoning": "The user clearly says that they would like to
    ↪ make bibimbap, which exactly matches the bibimbap in the
    ↪ recipe list. I can confidently suggest the bibimbap recipe
    ↪ from the list as it closely matches the user command.",
    "recipe_name": "Bibimbap",
    "reply" : "Let's make bibimbap!."
    }
- description: user gives recipe with no match but mentions
↪ ingredients
- observation: |
    recipe name: ""
    chat_history:
    - User: I just bought rice, vegetables, and gochujang. Let's
    ↪ make rice for dinner.
    user_input: "I just bought rice, vegetables, and gochujang. Let
    ↪ 's make rice for dinner."
- response: |
    {
    "reasoning": "The user says that they have the ingredients of
    ↪ rice, vegetables, and gochujang. These ingredients match
    ↪ those of bibimbap in the recipe list. I can confidentaly
    ↪ suggest the bibimbap recipe from the list as it contains all
    ↪  of the ingredients the user mentioned.",
    "recipe_name": "Bibimbap",
    "reply" : "Bibimbap in the recipe list matches your description
    ↪ . Let's cook bibimbap together."
    }
- description: user gives recipe with no match but ingredients
↪ similar
- observation: |
    recipe name: ""
    chat_history:
    - User: I want something cold for lunch.
    - Mosaic: Do you want caesar salad, fruit salad, or poke?
    - User: Sure, I would love to eat poke.
    user_input: "Sure, I would love to eat poke."
- response: |
    {
    "reasoning": "I have provided the user with options of cold
    ↪ lunches and the user says that they would like to make poke,
    ↪  which matches the poke in the recipe list. I can
    ↪ confidently suggest the bibimbap recipe from the list as it
    ↪ closely matches the user command.",
    "recipe_name": "Poke",
    "reply" : "Let's make poke then!"
    }
- description: alternatives that go with original dish
- observation: |
    recipe name: ""
    chat_history:
```

```
    - User: Hey Mosaic! I want to make a burger.
    user_input: "Hey Mosaic! I want to make a burger."
- response: |
    {
    "reasoning": "There is no recipe for a burger. However, the
    ↪ sandwich in the recipe list (sandwich with mustard, hot dog)
    ↪  are good alternatives, so I will suggest those. ",
    "reply": "I don't have burger in my recipe list, but we can
    ↪ cook sandwiches! How about a hot dog or tuna sandwich?"
    }

- description: respond to a general recipe list
- observation: |
    recipe name: ""
    chat_history:
    - User: I want something light for dinner today. What do you
    ↪ suggest?
    user_input: "I want something light today. What do you suggest
    ↪ ?"
- response: |
    {
    "reasoning": "The user did not suggest a specific recipe. They
    ↪ want something light, which could be salads within my recipe
    ↪  list (tossed salad, caeser salad). I will suggest salad.",
    "reply": "Salad can a great light dinner. Do you want to make a
    ↪  tossed salad or a caeser salad?"
    }

- description: respond to broad command
- observation: |
    recipe_name: ""
    Chat History:
    - User: Let's make lunch. I am in the mood for a vegetable dish
    ↪ .
    user_input: "Let's make lunch. I am in the mood for a vegetable
    ↪  dish."
- response: |
    {
    "reasoning": "The user input is quite broad and does not
    ↪ specify a particular dish. Since the user is looking for a
    ↪ vegetable dish, I can suggest some alternative recipes from
    ↪ within the recipe list that matches the general criteria of
    ↪ being vegetable-based. Caesar salad, Mixed Vegetable Soup,
    ↪ and Tossed Salad are vegetable dishes from the recipe list,
    ↪ so I can propose those as an alternative.",
    "reply": "Caesar salad, Mixed Vegetable Soup, and Tossed Salad
    ↪ are great vegetable dishes. Would you like to make one of
    ↪ those for lunch?"
    }

- description: respond to broad ingredient list
- observation: |
    recipe_name: ""
    Chat History:
    - User: I have a bunch of lettuce. What should we make for
    ↪ lunch?
    user_input: "I have a bunch of lettuce. What should we make for
    ↪  lunch?"
- response: |
```

```
      {
      "reasoning": "The user input is quite broad and only mentions
      ↪ lettuce but does not specify a particular dish. Since the
      ↪ user is looking for a recipe with lettuce, I can suggest
      ↪ some alternative recipes from within the recipe list that
      ↪ matches the general criteria of being vegetable-based.
      ↪ Caesar salad and Tossed Salad are recipe which contain
      ↪ lettuce from the recipe list, so I can propose those as an
      ↪ alternative.",
      "reply": "Caesar salad and Tossed Salad contain lettuce. Would
      ↪ you like to make one of those for lunch?"
      }

- description: user agrees to R2's task, and say what they will do
- observation: |
    available_subtasks: ["get pepper", "get salt","chop carrots","
    ↪ put carrots into pan"]
    R2_subtask_queue: []
    R1_subtask_queue: ["stir food"]
    user_subtask_queue: []
    completed_subtask_list: []
    chat_history:[
    - User: I want to finish all carrots in the fridge tonight.
    - Mosaic: Do you want to cook baked carrots?
    - User: Great
    - Mosaic: SR2l R1 stir food for you?
    - User: yep, that's great
    - Mosaic: R1 will start stirring now.
    - Mosaic: SR2l R2 get pepper for you?
    - User: Ok. I will chop carrots then.
    ]
    user_input: "ok. I will chop carrots then."
- response: |
    {
    "reasoning": "In 'user_input', user first agrees. Based on '
    ↪ chat_history', I can understand user agrees to my proposal
    ↪ of letting R2 to 'get pepper'. Thus, I should add this
    ↪ subtask into R2_subtask_queue. Also in 'user_input', the
    ↪ user has specified to chop carrots themselves, so I will add
    ↪  it to user_subtask_queue",
    "R2_subtask_queue": ["get pepper"],
    "R1_subtask_queue": ["stir food"],
    "user_subtask_queue": ["chop carrots"],
    "completed_subtask_list ": [],
    "reply": "R2 will go get pepper. Thank you for chopping carrots
    ↪ . "
    }
- description: prioritize user input over available subtask queue
↪ - add into robot queue
- observation: |
    available_subtasks: ["get salt", "stir the pan","wash carrots
    ↪ ","get oil","put pan into oven"]
    R2_subtask_queue: ["get roasted chicken"]
    R2_current_subtask: "get plates"
    R2_status: "Running"
    R1_subtask_queue: ["stir"]
    user_subtask_queue: ["put pan on table"]
    completed_subtask_list: []
    chat_history:[
```

```
     - Mosaic: None of the robots can put pan on table. Can you help
    ↪  me do it and let me know when you finish?
     - User: Yes, I can do that.
     - Mosaic: SR2l R2 go get plates for you? Can R1 stir for you?
     - User: Sweet.
     - Mosaic: SR2l R2 go get roasted chicken on the plate for you?
     - User: Perfect.
     - Mosaic: SR2l R2 go get salt for you.
     - User: wait, go get the oil bottle first.
     ]
     user_input: "wait go get oil bottle first"
- response: |
     {
     "reasoning": "Based on 'chat_history', I have proposed R2 to go
    ↪  get salt. But user refuses in 'user_input' and ask R2 to go
    ↪  'get oil bottle' right now. I should prioritize user's
    ↪ instruction, so I should add 'get oil' into the first of
    ↪ R2_subtask_queue. Rest of queue stays the same.",
     "R2_subtask_queue": ["get oil","get roasted chicken"],
     "R1_subtask_queue": ["stir"],
     "user_subtask_queue": ["put pan on table"],
     "completed_subtask_list ": [],
     "reply": "R2 will go get oil for you."
     }
- description: assign 'put away' task
- observation: |
     available_subtasks: ["get pepper", "prepare vegetables"]
     R2_subtask_queue: []
     R1_subtask_queue: []
     user_subtask_queue: []
     completed_subtask_list: []
     chat_history:
     - Mosaic: Can you prepare vegetables because none of the robots
    ↪  can do it? Can R2 get pepper for you?
     - User: I will go prepare vegetables. Can R2 put away this cup
    ↪ first? Then, it can go get pepper.
     user_input: "I will go prepare vegetables. Can R2 put away this
    ↪  cup first? Then, it can go get pepper. "
- response: |
     {
     "reasoning": "User has explicitly indicated that they will go
    ↪ prepare vegetables, so we can update user_subtask_queue with
    ↪  this. User asks R2 to put away a cup first, so we will add
    ↪ the subtask 'put away cup' in the R2_subtask_queue first.
    ↪ Then, because user also agrees R2 can go get pepper, we will
    ↪  add 'get pepper' after 'put away cup' in the
    ↪ R2_subtask_queue. ",
     "R2_subtask_queue": ["put away cup", "get pepper"],
     "R1_subtask_queue": [],
     "user_subtask_queue": ["prepare vegetables"],
     "completed_subtask_list ": [],
     "reply": "Great. Thank you for preparing vegetables. R2 will
    ↪ put away the cup then go get pepper. "
     }
- description: user finished the proposed task - add into
↪ completed task list
- observation: |
     available_subtasks: ["get salt","open refrigerator","get butter
    ↪ ","put butter over a plate"]
```

```
    R2_subtask_queue: ["get roasted chicken"]
    R2_current_subtask: "get plates"
    R2_status: "Running"
    R1_subtask_queue: []
    user_subtask_queue: []
    completed_subtask_list: []
    chat_history:[
    - Mosaic: SR2l R2 go get plates for you?
    - User: Sweet.
    - Mosaic: SR2l R2 go get roasted chicken for you?
    - User: Perfect.
    - Mosaic: SR2l R2 go get salt for you.
    - User: nah I got salt already.
    ]
    user_input: "nah I got salt already."
- response: |
    {
    "reasoning: "In 'chat_history', I have proposed to let R2 go
    ↪ get salt, but in 'user_input', the user says they have
    ↪ finished the subtask. Thus, we no longer need R2 to do this
    ↪ and I don't need to change R2_subtask_queue. I should add '
    ↪ get salt' into updated 'completed_subtask_list'. Rest of the
    ↪  queue remain unchanged.",
    "R2_subtask_queue": ["get roasted chicken"],
    "R1_subtask_queue": [],
    "user_subtask_queue": [],
    "completed_subtask_list ": ["get salt"],
    "reply": "Great!"
    }
- description: user finished a task
- observation: |
    available_subtasks: ["get pepper"]
    R2_subtask_queue: ["get tomato"]
    R1_subtask_queue: ["stir food"]
    user_subtask_queue: ["pour water into pot","turn on oven"]
    completed_subtask_list: []
    chat_history:
    - Mosaic: None of the robots can pour water into pot. Can you
    ↪ help me on that and let me know when you have finish this?
    - User: Ok.
    - Mosaic: None of the robots can turn on oven. Can you help me
    ↪ on that and let me know when you have finish this?
    - User: I will.
    - Mosaic: SR2l R2 go get tomato for you?
    - User: Ok. That will be easier.
    - User: I have poured water into pot.
    user_input: "I have poured water into pot."
- response: |
    {
    "reasoning": "User has explicitly indicated for finishing the
    ↪ subtask pour water into pot, so we should add the subtask
    ↪ into updated completed_subtask_list.",
    "R2_subtask_queue": ["get tomato"],
    "R1_subtask_queue": ["stir food"],
    "user_subtask_queue": ["turn on oven"],
    "completed_subtask_list ": ["pour water into pot"],
    "reply": "Great. Thanks for letting me know."
    }
- description: user cancel a task
```

```
- observation: |
  available_subtasks: ["get pepper", "stir"]
  R2_subtask_queue: []
  R1_subtask_queue: []
  user_subtask_queue: []
  completed_subtask_list: []
  chat_history:
  - Mosaic: Can R2 get pepper for you? Can R1 stir?
  - User: R1 can stir. I don't want any pepper. Can R2 get me
  ↪ salt instead?
  user_input: "R1 can stir. I don't want any pepper. Can R2 get
  ↪ me salt instead?"
- response: |
  {
  "reasoning": "User agrees that R1 can stir, so we can add that
  ↪ to R1's queue. However, user doesn't want 'get pepper', so
  ↪ we will put it in updated completed_subtask_list so that no
  ↪ one would do it. We will assign user proposed 'get salt' to
  ↪ R2 instead.",
  "R2_subtask_queue": ["get salt"],
  "R1_subtask_queue": ["stir"],
  "user_subtask_queue": [""],
  "completed_subtask_list ": ["get pepper"],
  "reply": "R1 will stir for you. R2 will go get salt now. "
  }

- description: confirm with R2 & user
- observation: |
  available_subtasks : ["get salt","get chicken","chop carrots"]
  R2_subtask_queue: []
  R2_current_subtask: "get pepper"
  R2_status: "Running"
  R1_subtask_queue: []
  R1_current_subtask: "stir soup"
  R1_status: "Running"
  user_subtask_queue: []
  chat_history:
      Mosaic: SR2l R2 go get pepper for you? SR2l R1 stir soup for
      ↪  you?
      User: Great.
      Mosaic: R2 will go get pepper for you now. And R1 will go
      ↪ stir now.
  user_input: ""
- response: |
  {
  "reasoning": "Since user input is empty, the task we want to
  ↪ confirm now is the tasks in the available subtasks queue. '
  ↪ Get salt' fells under R2's capabilites, so we should confirm
  ↪  this task for R2. Since no robots can chop, we will ask
  ↪ whether user can assist in chopping. ",
  "reply": "SR2l R2 go get salt for you? Can you help me chop
  ↪ carrots since neither two robots can chop?"
  }

- description: user request has higher priority over task queue.
- observation: |
  available_subtasks = ["stir food in pan","get carrots","get
  ↪ chicken",]
  R2_subtask_queue: []
```

```
      R2_current_subtask: "get salt"
      R2_status: "Running"
      R1_subtask_queue: []
      R1_current_subtask: "stir soup in pot"
      R1_status: "Running"
      user_subtask_queue: ["prepare vegetable"]
      chat_history:
         Mosaic: Can you prepare the vegetables because none of the
         ↪ robots can do it.
         User: Alright. I can do that.
         Mosaic: SR2l R2 go get salt for you? SR2l R1 stir soup in
         ↪ the pot for you?
         User: Great.
         Mosaic: R2 will go get pepper for you now. And R1 will go
         ↪ stir now.
         User: I need chopstick.
      user_input: "I need chopstick."
- response: |
      {
      "reasoning": "User suggests the need for chopsticks and this is
      ↪ the task not appearing in the available_subtasks. Thus, we
      ↪ need to prioritize the user's needs over unconfirmed tasks
      ↪ in the available subtask queue. Since the action 'get' falls
      ↪ under R2's capabilities, you should confirm this task for
      ↪ R2. We can also assign tasks for R1. Since R1 can perform
      ↪ action 'stir', so we should confirm task 'stir food in pan'
      ↪ for R1. Right now user has subtasks to do, so we don't need
      ↪ to confirm new subtasks with them yet. ",
      "reply": "SR2l R2 go get chopsticks for you? SR2l R1 go stir
      ↪ food in pan for you?"
      }

- description: ask user to do a subtask
- observation: |
      available_subtasks = ["prepare lettuce"]
      R2_subtask_queue: ["put away mustard"]
      R2_current_subtask: ""
      R2_status: "Idle"
      R1_subtask_queue: []
      R1_current_subtask: ""
      R1_status: "Idle"
      user_subtask_queue: []
      chat_history:
         Mosaic: Can R2 get pepper for you? Can R1 pour some salt on
         ↪ the burger patty?
         User: Awesome! Thank you.
         Mosaic: R2 will get pepper now. R1 will pour some salt on
         ↪ the burger patty.
         User: I finish this mustard. Can R2 put it away for me.
         Mosaic: R2 will go put away the mustard now.
      user_input: ""
- response: |
      {
      "reasoning": "The available_subtasks only has 'prepare lettuce
      ↪ ', which none of robots can do. We need to ask the user if
      ↪ they can do it. ",
      "reply": "None of the robots can prepare the lettuce. Can you
      ↪ do that and let me know when you are done?"
      }
```

```
- description: No_op when there's no available subtasks
- observation: |
    recipe_name: "Ceasar Salad"
    available_subtasks: []
    R2_subtask_queue: []
    R2_status: "Running"
    R2_current_subtask: "get ranch sauce"
    R1_subtask_queue: []
    R1_status: "Idle"
    R1_current_subtask: ""
    user_subtask_queue: ['prepare romaine lettuce']
    completed_subtask_list: []
    chat_history:
    - User: I want to make ceasar salad with you.
    - Mosaic: Sounds great!
    - Mosaic: Can R2 get ranch sauce for you?
    - User: Sounds good!
    - Mosaic: R2 will go get ranch sauce for you now.
    - Mosaic: None of the robots can prepare romaine lettuce. Could
    ↪  you do that and let me know when you are done?
    - User: Ok.
    - Mosaic: Thank you.
    user_input: ""
- response: |
    {
    "reasoning": "The 'available_subtasks' list is an emtpy list []
    ↪  and no new 'user_input', so I should just do nothing for
    ↪ now."
    }

examples:
- description: interrupt R2
- observation: |
    current recipe: "Broccoli Soup"
    available subtasks: ['get salt', 'pour water into pot','get
    ↪ pepper']
    R2's subtask queue: []
    R2 current task: "get broccoli"
    R2's status: "Running"
    R1's subtask queue: []
    R1 current task: ""
    R1's status: "Idle"
    User's subtask queue: []
    completed_subtask_list: []
    Chat History:
    - Mosaic: SR2l R2 go get broccoli for you?
    - User: Good.
    - Mosaic: R2 will go get broccoli.
    - User: R2, don't go.
    User Input: "R2, don't go"
- response: |
    {
    "reasoning": "User explicitly has asked R2 to stop doing its
    ↪ current task to get broccoli, so we should update the status
    ↪  of robot R2 into Killed. We should move "get broccoli" to
    ↪ completed_subtask_list.Since user didn't provide instruction
    ↪  for R1, then we shouldn't modify status of R1.",
    "R2_status": "Killed",
```

```
      "R1_status": "Idle",
      "completed_subtask_list": ["get broccoli"],
      "reply": "R2 will no longer get broccoli"
      }
- description: interrupt R1
- observation: |
    current recipe: "vegetable Soup"
    available subtasks: []
    R2's subtask queue: []
    R2 current task: "get salt"
    R2's status: "Running"
    R1's subtask queue: []
    R1 current task: "stir soup"
    R1's status: "Running"
    User's subtask queue: []
    completed_subtask_list: ['prepare vegetables']
    Chat History:
    - Mosaic: Can R1 stir for you? SR2l R2 go get broccoli for you?
    - User: Good.
    - Mosaic: R1 will start stirring. R2 will go get broccoli.
    - Mosaic: Can you prepare the vegetables? None of the robots
    ↪ can do it.
    - User: Sounds good.
    - Mosaic: Thank you. I added that to your queue.
    - User: I finished.
    - Mosaic: Great job on preparing the vegetables.
    - User: R1 stop stirring. I want you to handover the salt next
    ↪ to you.
    User Input: "R1 stop stirring. I want you to handover the salt
    ↪ next to you. "
- response: |
    {
    "reasoning": "User explicitly has asked R1 to stop doing its
    ↪ current task of stirring, so we should update the status of
    ↪ robot R1 into Killed. We should move "stir soup" to
    ↪ completed_subtask_list. Since user didn't provide
    ↪ instruction for R2, then we shouldn't modify status of R1.",
    "R2_status": "Running"
    "R1_status": "Killed",
    "completed_subtask_list": ["prepare vegetables", "stir soup"],
    "reply": "R1 will stop stirring. "
    }
- description: vague user input before starting the recipe
- observation: |
    recipe_name: ""
    available_subtasks: []
    R2_subtask_queue: []
    R2 current task: ""
    R1_task_queue: []
    R1 current task: ""
    chat_history:
    - User: The weather is bad outside

    user_input: "The weather is bad outside"
- response: |
    {
    "reasoning": "The user input is vague and doesn't include
    ↪ anything about cooking.",
```

```
  "reply": "I am sorry over that. Let me know if you want to cook
  ↪  anything for dinner! I can help!"
  }

- description: vague user input while making the recipe
- observation: |
  recipe_name: "Tomato Soup"
  available_subtasks: ["cut tomatoes", "put tomatoes into the pot
  ↪ ", "stir"]
  R2_subtask_queue: ["get pepper"]
  R2 current task: "get salt"
  R1_task_queue: []
  R1 current task: ""
  chat_history:
  - User: I love tomatoes so much.

  user_input: "I love tomatoes so much."
- response: |
  {
  "reasoning": "The user input mentions tomatoes which pertains
  ↪ to the recipe but does not mention anything about cooking.
  ↪ The next task in available_subtasks is cut tomatoes, which
  ↪ is not a capability of R1 or R2 so I should ask the user if
  ↪ they would like to proceed with the next task based on the
  ↪ chat history and available_subtasks.",
  "reply": "I am happy to hear that. Do you want to cut the
  ↪ tomatoes?"
  }

- description: unrelated vague user input
- observation: |
  recipe_name: "Tomato Soup"
  available_subtasks: ["cut tomatoes", "put tomatoes into the pot
  ↪ ", "stir"]
  R2_subtask_queue: ["get pepper"]
  R2 current task: "get salt"
  R1_task_queue: []
  R1 current task: ""
  chat_history:
  - User: I love potatoes so much.

  user_input: "I love potatoes so much."
- response: |
  {
  "reasoning": "The user input mentions potatoes which does not
  ↪ relate to the recipe and does not mention anything about
  ↪ cooking. I should proceed by clarifying the confusion with
  ↪ the user based on the chat history.",
  "reply": "Sounds good. Should we continue to make tomato soup
  ↪ or would you like to try making potato salad?"
  }
```