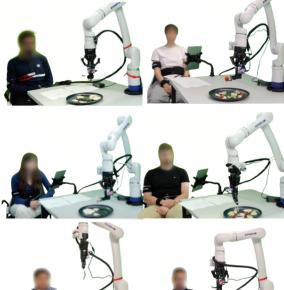


# A Human-in-the-loop Confidence-Aware Failure Recovery Framework for Modular Robot Policies

Anonymous Authors

## In Lab Studies - Individuals With Emulated Mobility Limitations



## In Home Studies - Individuals With Mobility Limitations

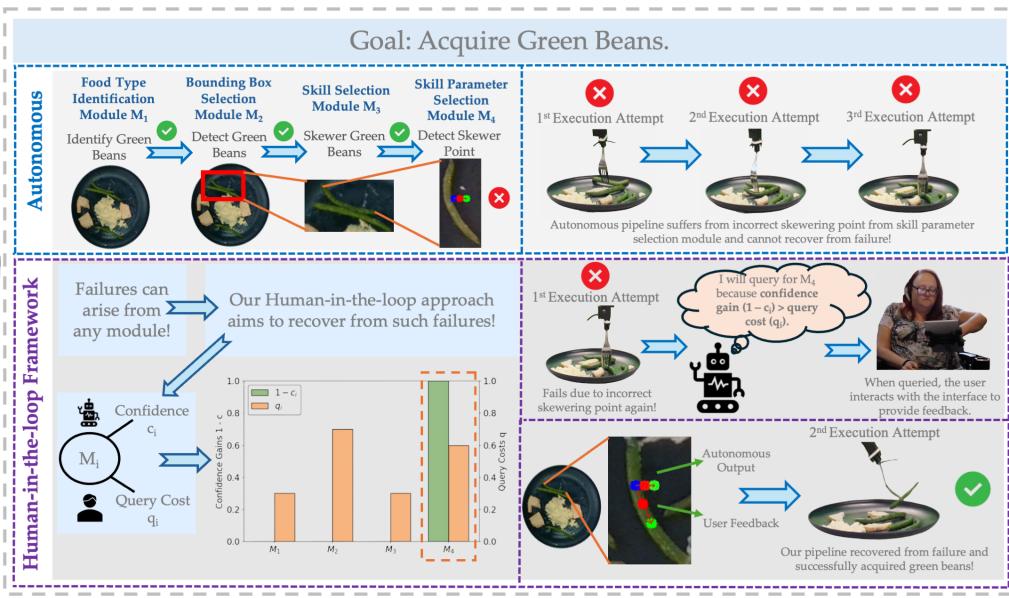


Figure 1: Our human-in-the-loop framework for failure recovery leverages confidence estimates from a modular policy, along with predicted estimates of user workload, to decide what to ask and when to act autonomously.

## Abstract

Robots operating in unstructured human environments inevitably encounter failures, especially in robot caregiving scenarios. While humans can often help them recover, excessive or poorly targeted queries impose unnecessary cognitive and physical workload on the human partner. We present a human-in-the-loop failure-recovery framework for modular robotic policies, where a policy is composed of distinct modules such as perception, planning, and control, any of which may fail and often require different forms of human feedback. Our framework integrates calibrated estimates of module-level uncertainty with models of human intervention cost to decide which module to query and when to query the human. It separates these two decisions: a module selector identifies the module most likely responsible for failure, and a querying algorithm determines whether to solicit human input or act autonomously.

We evaluate several module-selection strategies and querying algorithms in controlled synthetic experiments, revealing trade-offs between recovery efficiency, robustness, and interaction workload. Finally, we deploy the framework on a robot-assisted bite acquisition system and demonstrate, in studies involving individuals with both emulated and real mobility limitations, that it improves recovery success while reducing the workload imposed on users. Our results highlight how explicitly reasoning about both robot uncertainty and human effort can enable more efficient and user-centered failure recovery in collaborative robots. See our website: <https://biteacquisitionrobot6-art.github.io/>

## ACM Reference Format:

Anonymous Authors. 2026. A Human-in-the-loop Confidence-Aware Failure Recovery Framework for Modular Robot Policies. In *Proceedings of ACM/IEEE International Conference on Human-Robot Interaction (HRI) (HRI '26)*. ACM, New York, NY, USA, 18 pages. <https://doi.org/XXXXXX.XXXXXXXX>

Unpublished working draft. Not for distribution. This version is intended for individual research and study only. It may not reflect final published versions of the work. Copying or further distribution of this version is expressly prohibited without the permission of the copyright holders.

HRI '26, Edinburgh, Scotland, UK

© 2026 Copyright held by the owner/author(s). Publication rights licensed to ACM.

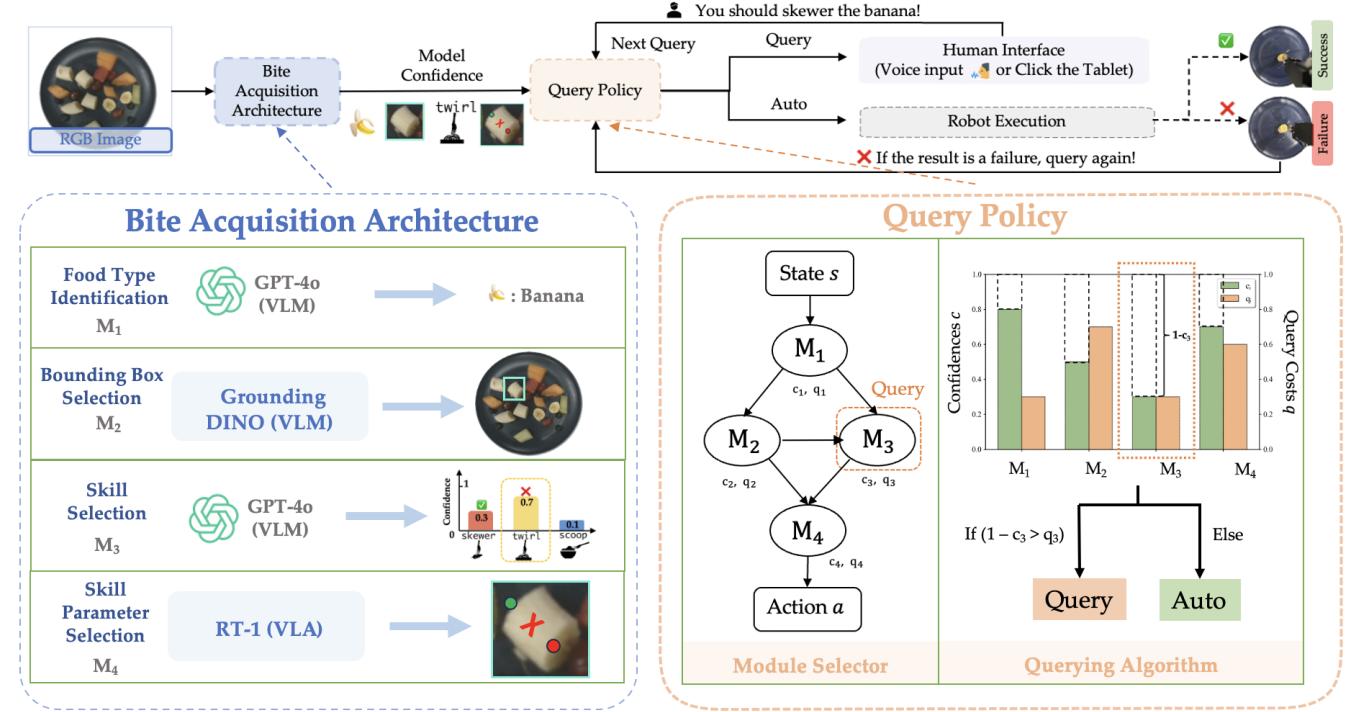
ACM ISBN 978-1-4503-XXXX-X/2018/06

<https://doi.org/XXXXXX.XXXXXXXX>

2026-01-02 02:45. Page 1 of 1-18.

## 1 Introduction

Robots deployed in the wild inevitably fail, especially when assisting individuals with mobility limitations in performing activities of daily living (ADLs) in unstructured home environments [13, 36, 39]. While humans can often help them recover [3, 10], requesting help too frequently imposes cognitive and physical workload on the human partner [16, 35]. Designing robots that know when and what



**Figure 2: Overall human-in-the-loop decision failure recovery framework, grounded in the robot-assisted bite acquisition domain. The recovery framework first invokes a module selector to decide which of the modules to query for (e.g the bounding box selector or the skill parameter selector. The framework then invokes a querying algorithm, which decides whether to ask the user for help or act autonomously.**

to ask humans is therefore central to effective human–robot interaction. In this work, we focus on modular robot systems, rather than end-to-end models such as vision–language–action (VLA) policies, because modular architectures are typically more interpretable and thus more amenable to structured failure recovery. However, failures in modular policies can arise in any of the perception, planning, or control modules, each of which may require a different form of human feedback, making it challenging to effectively incorporate human input during recovery.

Recovering from failures in such modular systems is difficult for two reasons. First, identifying which module has failed is non-trivial; an error in perception can cascade into planning and execution, making it unclear where intervention is most effective. Second, the confidence scores provided by individual modules are often miscalibrated and do not reliably predict whether a task will succeed. A robot that queries for help too frequently risks frustrating the human partner and increasing interaction costs, while one that queries too little risks repeated failures that undermine both task performance and trust. Balancing this trade-off between task success and human workload lies at the heart of collaborative autonomy. By using *confidence-aware reasoning* where the robot uses its own estimates of uncertainty about its decisions together with models of the workload on the human [3, 11, 31], we can design *workload-aware interaction* strategies that recover from failures

more efficiently while maintaining high user satisfaction. This capability is particularly crucial in assistive settings, such as in robot caregiving, where both reliable task performance and the quality of the user’s experience determine the acceptability and long-term use of robot assistance.

To address this challenge, we propose a human-in-the-loop failure-recovery framework for modular robot policies that integrates calibrated estimates of module-level uncertainty with models of human intervention cost. The framework makes two key decisions at every recovery attempt. First, a *module selector* determines which component of the modular policy to query, e.g. whether to seek human correction for a perception error or for a planning decision. This choice is necessary because failures may arise in multiple modules, and querying all of them would be inefficient and would impose unnecessary workload on the human partner. Second, a *querying algorithm* determines whether to query the human at all or to act autonomously at a given point. This decision is critical for trading off the expected gain in task-level success against the interaction cost to the human partner.

We formulate these decisions within a sequential decision-making framework, in which the robot weighs anticipated reductions in task-level uncertainty against the expected cost of additional human interaction. We evaluate several module-selection strategies, including brute-force and graph-based methods that exploit the structure and redundancy of the modular policy, and we compare

alternative querying algorithms that govern when to solicit input to balance recovery efficiency, robustness, and workload. Finally, we deploy the framework on a robot-assisted bite acquisition system [3, 13, 14, 36], demonstrating improved recovery success and reduced user querying workload across studies involving individuals with both emulated and real mobility limitations.

Our work makes the following novel contributions:

- We propose a human-in-the-loop failure-recovery framework for modular robotic policies that explicitly integrates the interaction aspect of recovery. The framework makes two complementary decisions: a *module selector* chooses which component of the policy to query, and a *querying algorithm* decides whether and when to solicit human input versus act autonomously.
- We incorporate calibrated estimates of module-level uncertainty together with models of human intervention cost to guide both decisions, enabling robots to recover more efficiently and to avoid unnecessary queries.
- We conduct systematic evaluations of alternative module-selection strategies (e.g., brute-force, graph-based, binary-tree, and mixed-integer programming (MIP) selectors) and querying algorithms (e.g., execute-first, query-then-execute, query-until-confident, and workload-aware variants) in controlled synthetic experiments, revealing trade-offs in robustness, efficiency, and interaction workload.
- We demonstrate the approach in a real-world robot-assisted bite acquisition task, achieving higher recovery success and lower user querying workload in both in-lab studies with emulated limitations and an in-home study involving users with real mobility limitations.

## 2 Related Work

**Anomaly detection and recovery in robotics.** Robotic policies, whether end-to-end or modular, are prone to failure during execution, especially in unstructured environments or when encountering novel objects. Despite recent advances in end-to-end robot learning, these approaches can be difficult to interpret [1, 8, 40]. In contrast, modular methods can offer more predictable and controllable performance in certain scenarios. In this work, we focus on full-stack modular methods that consider failure detection and recovery across different modules of a robot system. It is natural to represent such problems using a graph representation, as in prior work on online approaches for robot failure detection and recovery [2, 7, 28]. Additionally, the overall interaction process can be modeled by certain types of graphs, ranging from finite-state machines [7], MDPs/POMDPs [12, 28], and Behavior Trees [2, 6, 25]. Although graph representations have been widely studied in prior work, our method goes further by not only modeling systems using graphs but also leveraging graph algorithms to select which modules to query for information or intervention.

In general, there are two primary gaps in prior research that our approach addresses. First, while some prior work [2] assumes that the system is in a particular failure state, our approach includes failure detection through the use of confidence estimates. Second, most prior work focuses on autonomous recovery strategies, whereas our work develops human-in-the-loop approaches. We developed four

methods to model the modular autonomous policy and to select modules to be asked. Two of the four methods are graph-based methods.

**Robots that ask for help.** Robot systems may decide to query humans on the basis of several different criteria. Some approaches ask for help when a failure is detected after a robot action has been executed [21, 37], some based on a confidence estimate about potential robot actions to take [26, 32], and others based on an estimate about the potential information gain that could arise from soliciting human feedback [10, 29]. Departing from existing work, our query rules depend on a combination of both model confidences of different modules and estimated workload costs of querying.

The types of feedback the robot can ask for can also vary, ranging from natural language [5, 32, 34], to actions [21, 37, 38], to labels [5]. Our framework that detects the module to be queried is agnostic to the specific types of feedback. Unlike prior work, we simultaneously consider asking for help for multiple modules in a full-stack policy architecture, not just for a single action-selection module [19, 32] or reward model [10].

## 3 Problem Formulation

**Module graph.** We assume that our policy architecture  $\mathcal{P} : S \rightarrow A$  is structured in a module graph  $G_M$ , which describes the relationships between the  $N$  modules in the architecture. Here,  $S$  represents the policy input state (e.g. RGB-D images), and  $A$  represents the policy action output (e.g. robot end-effector trajectories). The vertices in  $G_M$  are modules, and an edge exists from one module  $M_i$  to another module  $M_j$  if the output of  $M_i$  is an input to  $M_j$ .

A *module* is a tuple  $M = (\pi, \mathcal{X}, \mathcal{Y}, c, q(t))$  where  $\pi : \mathcal{X} \rightarrow \mathcal{Y}$  is a policy function,  $c : \mathcal{X} \rightarrow [0, 1]$  is a confidence score, and  $q(t)$  is a time-dependent cost for querying the human about the module. For instance, a module could be a perception module that extracts geometric information about the robot's environment from RGB-D input, or a planning module that extracts a robot trajectory given obstacle states in the environment. We assume that one module in  $G_M$  (which we denote  $M_1$ ) takes in the state  $s \in S$  (i.e.  $\mathcal{X} = S$ ), and one module (which we denote  $M_N$ ) produces a final action  $a \in A$  (i.e.  $\mathcal{Y} = A$ ).

**Research problem.** Our goal is to decide when and what to query, which we decompose into (1) a module selector algorithm  $\psi_{ms}$  which chooses a single module, and (2) a querying algorithm  $\psi_q$  that decides whether to query or to execute the autonomous action  $a = \mathcal{P}(s)$ . For a particular state  $s \in S$ , the module selector algorithm  $\psi_{ms}$  and querying algorithm  $\psi_q$  should minimize  $J(\psi_{ms}, \psi_q)$ , the weighted sum of the querying cost  $J_{query}(t; \psi_{ms}, \psi_q)$  and the task cost  $J_{task}(t; \psi_{ms}, \psi_q)$ , over the interaction horizon  $T$ :

$$\begin{aligned} J(\psi_{ms}, \psi_q) &= \frac{1}{T} \sum_{t=1}^T (w J_{query} + (1-w) J_{task}) \\ &= \frac{1}{T} \sum_{t=1}^T \left( w \sum_{M_i \in \Omega(t)} q_i(t) + (1-w)(1-r_t) \right) \end{aligned} \quad (1)$$

where  $w \in [0, 1]$  determines the trade-off between minimizing human workload and recovering from failures efficiently,  $\Omega(t)$  denotes the set of modules for which we query at time  $t$  (and  $M_i$

349 refers to an individual module in this set),  $q_i(t)$  is the query cost  
 350 for module  $i$ , and  $r(t)$  is the task reward, which is equal to 1 if the  
 351 execution succeeds, and 0 otherwise.

352 **Human feedback.** We assume that for every module  $M_i$ , the  
 353 user can provide feedback  $f \in \mathcal{Y}_i$  by choosing a possible output of  
 354 that module (e.g. classifying an object or selecting a high level skill).  
 355 When we receive this feedback, we replace the module output with  
 356 the expert feedback  $f$ , along with an expert confidence  $c = c_{expert}$ .

357 **Module redundancy.** The overall task reward  $r_t$  depends on  
 358 the success of individual modules, which can combine in two ways:

- A *redundant* manner, where at least one module  $M_i$  must succeed for their combination to succeed, so success is a union. For example, a policy architecture may include two redundant modules that predict object material properties: an off-the-shelf foundation model, and a domain-specific neural network [39].
- A *non-redundant* manner, where all modules  $M_i$  must succeed for their combination to succeed, so success is an intersection. For example, a goal-reaching manipulation policy architecture may require two modules to succeed: a goal inference module, and a planning module that infers a trajectory to the goal [24].

368 The overall success  $r(t)$  can depend on both redundant and non-  
 370 redundant combinations of modules.

## 4 Module Selectors

372 The module selector algorithms  $\psi_{ms}$  select which module to query  
 373 by considering the 1-timestep version of the overall objective (Eq.  
 374 1). While we cannot determine *a priori* whether each module  $M_i$   
 375 is in failure or success state, we have access to confidences  $c_i$ . We  
 376 consider three *proxy objectives* to minimize:

377 *Proxy objective 1* (product-of-confidences):

$$w \sum_{M_i \in \Omega(t)} q_i(t) + (1 - w)(1 - \prod_{M_i \notin \Omega(t)} c_i) \quad (2)$$

378 *Proxy objective 2* (sum-of-uncertainties):

$$w \sum_{M_i \in \Omega(t)} q_i(t) + (1 - w) \sum_{M_i \notin \Omega(t)} (1 - c_i) \quad (3)$$

384 *Proxy objective 3* (redundancy-dependent):

$$w \sum_{M_i \in \Omega(t)} q_i(t) + (1 - w)(1 - \hat{r}(\Omega(t))) \quad (4)$$

387 where reward estimate  $\hat{r}(\Omega(t))$  combines sums (redundant) and  
 388 products (non-redundant) of confidences for modules  $M \notin \Omega(t)$   
 389 (see Appendix for justifications for how each proxy objective ap-  
 390 proximates  $r_{task}$ ).

391 We consider four module selector algorithms—two direct proxy  
 392 optimization algorithms and two graph based algorithms—and three  
 393 additional baseline module selectors, two of which do not consider  
 394 query costs or confidence scores, and one that only considers confi-  
 395 dence scores. The baselines are (1) **Never Query** which does not  
 396 choose any module to query, (2) **Topo Query** which simply selects  
 397 the first module in the topological ordering of the module graph  
 398  $G_M$  which has not been previously selected, and (3) **Confidence**  
 399 **Query** which selects the module with the lowest confidence score.

### 4.1 Proxy Objective Optimization

400 The **Mixed-Integer Programming (MIP)** module selector directly  
 401 minimizes the redundancy-dependent proxy objective (Eq. 4) using  
 402 a mixed-integer solver (SCIP [4]). The optimization variable is an  
 403 indicator on which module to query.

407 The **Brute-Force** module selector also directly evaluates the  
 408 redundancy-dependent proxy objective (Eq. 4) for querying each  
 409 module  $M_i$  in the module graph  $G_M$ . Then, it selects the module  
 410  $M_i$  that achieves the minimal objective.

### 4.2 Binary Tree Query

414 This algorithm uses a binary tree to model the yes/no decision of  
 415 querying each module. The edge weights are defined so that paths  
 416 from the root to the leaves correspond to the product-of-confidences  
 417 proxy objective (Eq. 2). We then treat the module optimization  
 418 problem as a shortest-cost path problem, running Dijkstra's on the  
 419 graph to obtain the module  $M_i$ .

420 The binary tree  $B = (V, E)$  has edges  $E$  corresponding to query-  
 421 ing or acting autonomously for module  $M_i$ , and vertices  $V$  for  
 422 query histories prior to module  $M_i$ . The edge costs are  $C(e) = q_i(t)$   
 423 for query actions and  $C(e) = -c_i \prod_{j \in prev} c_j + \prod_{j \in prev} c_j$  for au-  
 424 tonomous actions, where  $prev$  are prior autonomous modules on  
 425 the path from root to  $M_i$  (if  $prev$  is empty,  $C(e) = 1 - c_i$ ). Fig. 3  
 426 (left) shows an example with  $N = 3$ .

### 4.3 GraphQuery

428 This algorithm constructs a directed graph that represents the mod-  
 429 ular policy and treats module selection as a shortest path problem,  
 430 running Dijkstra's on the graph to select the module  $M_i$ .

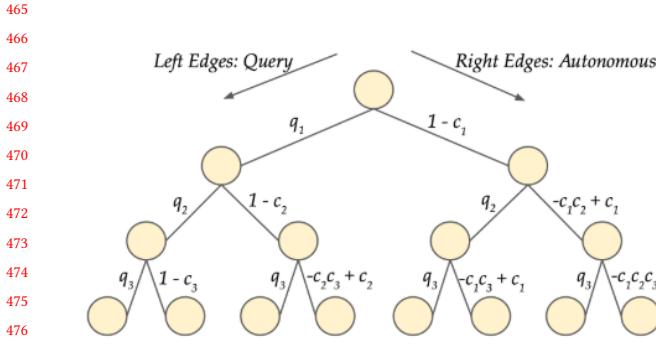
433 Given an ordering of modules in the module graph  $G_M$ , the graph  
 434  $G = (V, E)$  is shown in Fig. 3 (middle), where vertices  $V$  correspond  
 435 to the state of each module, and edges  $E$  correspond to actions,  
 436 both query  $a_{query}$  and autonomous  $a_{auto}$ . The vertices  $V$  consist  
 437 of an initial state  $s_{init}$ , pairs of success and failure states for each  
 438 module  $M_i$  ( $s_{M_i,success}$  and  $s_{M_i,failure}$ ). The edge costs are equal  
 439 to a scaled version of the query cost  $C(e) = \epsilon q_i(t)$  for edges  $e \in E$   
 440 corresponding to query actions. For edges  $e \in E$  corresponding to  
 441 autonomous actions, we set the cost equal to the uncertainty  $1 - c_i$ .  
 442 This method is equivalent to optimizing the sum-of-uncertainties  
 443 objective (Eq. 3).

## 5 Querying Algorithms

446 Querying algorithms decide whether to make a query or execute  
 447 the action  $a$  produced by the final action module  $M_N$  (Sec. 3). We  
 448 consider four querying variants illustrated in Fig. 3 (right), and one  
 449 baseline. They use the following primitives:

- **FORWARDPASS:** represents passing the current state  $s \in \mathcal{S}$  and  
 451 query set  $\Omega(t)$  through all of the modules in the module graph  
 452  $G_M$  to yield the final action  $a$  and updated confidences. In the  
 453 manipulation setting, this could represent providing the current  
 454 perceptual state (e.g. RGB-D camera data) to the policy archi-  
 455 tecture, producing the end-effector control sequence along with  
 456 intermediate module confidences.
- **EXECUTE:** represents executing the action  $a$  and observing overall  
 458 boolean *outcome*. In the manipulation setting, this could repre-  
 459 sent executing the end-effector control sequence on the robot  
 460 and assessing whether the manipulation task was successful.

462 **Execute-First** initially executes the autonomous action  $a$  from  
 463 the policy architecture  $\mathcal{P}$ . If  $a$  fails, it alternates between calling



**Figure 3:** (left) **BinaryTreeQuery** graph example for  $N = 3$ , (middle) **GraphQuery** graph example for  $N = 4$ , (right) **Querying algorithms**, which decide when to ask and when to execute actions, including **EXECUTE-FIRST**, which executes once prior to starting to query, **QUERY-THEN-EXECUTE**, which alternates between querying and execution, and **QUERY-UNTIL-CONFIDENT/QUERY-UNTIL-CONFIDENT-WORKLOAD-AWARE**, both of which repeatedly query until a stopping condition is met.

the module selector  $\psi_{ms}$  and executing until success. **Query-then-Execute** always alternates between invoking the module selector  $\psi_{ms}$  and executing until success.

#### Algorithm 1 EXECUTE-FIRST and QUERY-THEN-EXECUTE.

```

1: if EXECUTE-FIRST then
2:   a = FORWARDPASS(∅)
3:   outcome = EXECUTE(a)
4: else if QUERY-THEN-EXECUTE then
5:   outcome = FALSE
6: end if
7: while outcome = FALSE do
8:   Ω(t) = ψms(GM)
9:   a = FORWARDPASS(Ω(t))
10:  outcome = EXECUTE(a)
11: end while

```

**Query-until-Confident** repeatedly invokes  $\psi_{ms}$  until the proxy task objective (Eq. 4) exceeds a certain threshold  $\tau$ , then executes the resultant action  $a$ . **Query-until-Confident-Workload-Aware** repeatedly invokes  $\psi_{ms}$  until the confidence gain due to querying module  $M_i$  ( $c_{expert} - c_i$ ) is less than the scaled cost of querying  $\epsilon q_i(t)$ . The **Query-For-All** baseline repeatedly calls  $\psi_{ms}$  until it has queried for all modules, then executes the resultant action  $a$ .

## 6 Synthetic Simulation: Systematic Ablations

We develop a synthetic module simulation to investigate how system and user variables affect module selector and querying algorithm performance (Sec. 4, 5). Our simulation uses  $N$  modules connected via a random module graph  $G_M$ , where modules are logic gates (AND/OR) operating on Boolean inputs (including a Boolean state  $s$ ). Individual modules are set to be in either a success or failure state: modules in success states output their logic gate value, while modules in failure states always output FALSE. See Appendix F.1 for hyperparameter details (including proxy objective weighting and GraphQuery  $\epsilon$ ).

---

#### Algorithm 2 QUERY-UNTIL-CONFIDENT and QUERY-UNTIL-CONFIDENT-WORKLOAD-AWARE.

---

```

1: outcome = FALSE
2: while outcome = FALSE do
3:   repeat
4:     Ω(t) = ψms(GM)
5:     if QUERY-UNTIL-CONFIDENT then
6:       condition = [̂r(Ω(t)) > τ]
7:     else if QUERY-UNTIL-CONFIDENT-WORKLOAD-AWARE then
8:       condition = [cexpert - ci < εqi(t)]
9:     end if
10:    until condition
11:    a = FORWARDPASS(Ω(t))
12:    outcome = EXECUTE(a)
13: end while

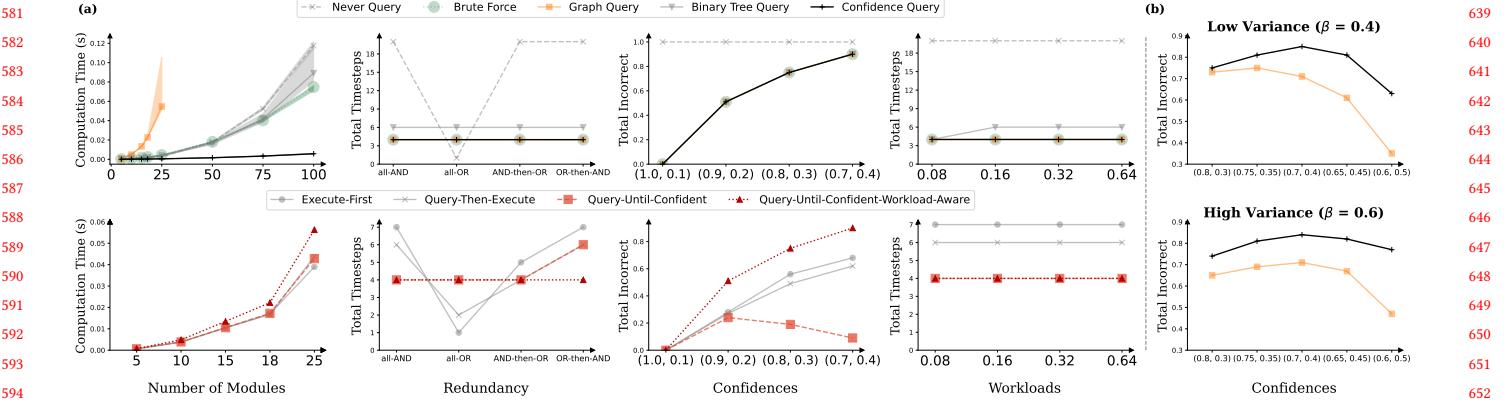
```

---

**Independent variables.** We examine four variables covering key system and user factors in human-robot interaction (full settings in Appendix F.1). We consider three system variables:

- **The number of modules  $N$ .** The number of modules in a robot policy architecture can vary, ranging from end-to-end architectures with a single module (e.g. a single VLA model) to more complex modular architectures.
- **Module redundancy/graph structures  $G_M$ .** Robot policy architectures include both redundant and non-redundant structures (Sec. 3). Different module selectors have different redundancy assumptions based on their proxy objective. We consider 4 redundancy structures with varied module types and orderings: fully redundant (all-OR), fully dependent (all-AND), fully redundant followed by fully dependent (OR-then-AND), fully dependent followed by fully redundant (AND-then-OR).
- **Confidence values  $c_i$ .** Robot policy architectures can have a diversity of confidence scores across the modules. We assign either a high or a low confidence value to each module in  $G_M$ . We additionally treat a module's confidence value as the probability that the module is in a success state.

Published in HRI '26, March 16–19, 2026, Edinburgh, Scotland, UK. Copyright © 2026, the author(s). Licensee: Association for the Advancement of Artificial Intelligence. This article is an open-access publication.



**Figure 4: (a) Systematic ablation experiments, showing the 4 independent variables in our simulations: (1) number of modules, (2) graph redundancies, (3) confidences, (4) query costs. We find that the BRUTEFORCE, GRAPHQUERY, and CONFIDENCEQUERY module selectors are the most robust to varying redundancy, confidences, and query costs, with CONFIDENCEQUERY being the most scalable. Additionally, we find that the QUERY-UNTIL-CONFIDENT and QUERY-UNTIL-CONFIDENT-WORKLOAD-AWARE querying algorithms are the most robust across redundancy and query costs, with QUERY-UNTIL-CONFIDENT having the best scalability and robustness to confidences; (b) Module heterogeneity experiments. We find that GRAPHQUERY outperforms CONFIDENCEQUERY, particularly when module confidences overlap and workload variance  $\beta$  is high.**

And we consider one user variable:

- **Query costs  $q_i$ .** The costs of querying for each module in a policy architecture can vary, due to the diversity of feedback types and the user's initial workload. We consider different levels of uniform query costs  $q_i$  for all modules.

We vary each independent variable in isolation, fixing the other variables as follows: (1)  $N = 10$  modules, (2) Fully dependent (all-AND) module structure, (3) 3 modules with confidence 0.1 and all other modules with confidence 1, (4) Uniform query cost of 0.32. We also assume no noise in the expert ( $c_{expert} = 1.0$ ). Additionally, we use a fixed querying algorithm (QUERY-UNTIL-CONFIDENT-WORKLOAD-AWARE) when analyzing module selectors, and a fixed module selector (GRAPHQUERY) when analyzing querying algorithms.

**Metrics.** We evaluate the methods using 5 metrics (lower is better):

- **Query Cost.**  $\sum_{t=1}^T \mathbb{1}\{\text{queried for } M_i \text{ at } t\} q_i(t)$ , where  $T$  is the number of timesteps needed to achieve success or the maximum time horizon (proportional to  $N$ ). This represents the total user workload required to correct the robot's failure.
- **Number of Failed Attempts.**  $\sum_{t=1}^T \mathbb{1}\{\text{failed at time } t\}$ . This represents the number of failed robot executions required to recover from the failure.
- **Computation Time.**  $\sum_{t=1}^T (\text{Time}(t))$ , where  $\text{Time}(t)$  is the total runtime between successive EXECUTE statements. This represents the total computation time required by the robot to run the module selector and the querying algorithm at each timestep (not including the time required to query).
- **Total Incorrect.** 0 if the agent successfully recovered from the failure after  $T$  timesteps, 1 otherwise. This represents whether the robot was able to recover from the failure given its time horizon constraint.

- **Total Timesteps  $T$ .** This includes both the number of queries and the number of failed executions required for the robot to recover from the failure.

## 6.1 Varying number of modules $N$

We find that CONFIDENCEQUERY is the most scalable module selector with linear time complexity in  $N$  (Fig. 4(a)), making it the most computationally efficient for large robot policy architectures.<sup>1</sup> As Computation Time scales with  $T$ , methods such as NEVERQUERY that require many timesteps incur higher computation time.

We note that QUERY-UNTIL-CONFIDENT-WORKLOAD-AWARE has the highest Computation Time (Fig. 4(a)), due to the additional step to predict the cost of querying, while the other three querying algorithms have nearly-identical values for Computation Time. Nevertheless, all querying algorithms have the same Total Query Cost (0.96) and Total Incorrect (0), regardless of  $N$ . For large robot policy architectures, it is thus advantageous to not select QUERY-UNTIL-CONFIDENT-WORKLOAD-AWARE if runtime is crucial.

## 6.2 Varying graph structure $G_M$

We find that all metrics are lower for the all-OR structure (across module selectors and querying algorithms). For this redundant structure, there is little benefit to ask for help to recover from failures, leading NEVERQUERY and EXECUTE-FIRST to perform the best. Module selector rankings (Fig. 4(a)) remain invariant to redundancy structure. We find that GRAPHQUERY is competitive with the other module selectors, even when its proxy objective does not match the redundancy structure.

<sup>1</sup>Because the computation time of the GRAPHQUERY exceeded 1 second for graph sizes  $N \geq 50$ , we omit the computation time for this module selector. Additionally, we omit reporting MIP results for all graph sizes because its computation time also exceeds this threshold.

We find that `QUERY-UNTIL-CONFIDENT-WORKLOAD-AWARE` and `QUERY-UNTIL-CONFIDENT` are generally the most efficient, except in the all-OR setting (where `EXECUTE-FIRST` and `QUERY-THEN-EXECUTE` outperform them). Thus, if the robot policy architecture is not fully redundant, we recommend either `QUERY-UNTIL-CONFIDENT` variant to maximize recovery efficiency.

### 6.3 Varying confidence values $c_i$

We find that all metrics (e.g. *Total Incorrect*, Fig. 4(a)) are higher when confidences overlap (across module selectors and querying algorithms). As `BRUTEFORCE`, `GRAPHQUERY`, and `CONFIDENCEQUERY` are the most competitive module selectors regardless of confidence level, we recommend any of these module selectors.

We find that the `QUERY-UNTIL-CONFIDENT` querying algorithm is the most robust to confidence variations. We recommend any querying algorithm if the confidence scores are well-separated in the robot policy architecture, and `QUERY-UNTIL-CONFIDENT` if confidence values are overlapping.

### 6.4 Varying query costs $q_i$

We find that all metrics except *Total Incorrect* are invariant to  $q_i$ , with `BRUTEFORCE`, `GRAPHQUERY`, and `CONFIDENCEQUERY` having lower *Total Timesteps* (Fig. 4(a)) compared to the other module selectors. Regardless of the level of user querying workload, we recommend either of these module selectors.

We find that besides *Total Query Cost*, all metrics are largely invariant to  $q_i$  across querying algorithm. As both `QUERY-UNTIL-CONFIDENT` variants have the lowest *Total Failed Attempts* and generally lower *Total Timesteps* (Fig. 4(a)), we recommend either of these variants.

## 7 Synthetic Simulation: Module Heterogeneity

We now consider a simulation setting with an additional user variable, **query cost variance**  $\beta$ , allowing the query costs  $q_i$  to vary across modules. This models more realistic cost variation across feedback types. We sample each  $q_i$  uniformly around a nominal value  $Q$ , i.e.,  $q_i \sim \text{Uniform}[(1 - \beta)Q, (1 + \beta)Q]$ , with  $Q = 0.32$ . We compare the two best performing module selectors from Sec. 6: the query cost-aware method `GRAPHQUERY`, and the confidence-only method `CONFIDENCEQUERY` (with the `QUERY-UNTIL-CONFIDENT-WORKLOAD-AWARE` querying algorithm).

We find that `GRAPHQUERY` consistently outperforms `CONFIDENCEQUERY` in the *Total Timesteps* and *Total Incorrect* metrics, especially when module confidences overlap or workload variance is high (*Total Incorrect* in Fig. 4(b)). Under these conditions, upstream modules with lower confidence and high workload may succeed, whereas downstream modules with higher confidence but low workload may fail. `CONFIDENCEQUERY` incorrectly selects these upstream modules first, whereas `GRAPHQUERY` correctly selects the downstream module (full results in Appendix G).

## 8 Robot-Assisted Bite Acquisition Experiments

Our robot-assisted bite acquisition architecture consists of  $N = 4$  modules: food type identification (GPT-4o), bounding box selection (GroundingDINO), skill selection (GPT-4o), and skill parameter selection (RT-1). The first three modules are VLMs, while the skill

selection module is a VLA. Food type identification and skill selection support learning from feedback via retrieval-augmented generation (RAG) (architecture and RAG details in Appendix C). We estimate module query costs using a predictive workload model from prior work [3] (Appendix E).

To select a module selector and query algorithm for bite acquisition, we adopt the recommendation from the closest synthetic setting as follows. We use the  $N = 10$  setting to approximate system scale; we model bite acquisition as non-redundant (all-AND), since all modules must be correct for success (Sec. 3); we use the confidence value setting (1, 0.1) to match our binary calibrated module confidence scores (Appendix C.2); we select  $q_i = 0.32$  to match the empirical mean of workload model predictions; we assume  $c_{expert} = 1.0$  as expert feedback is available for all modules. Based on these conditions (Sec. 6), along with the query cost variance experiments (Sec. 7), we use the `GRAPHQUERY` module selector with the `QUERY-UNTIL-CONFIDENT-WORKLOAD-AWARE` querying algorithm, which performed best in *Total Timesteps* and *Total Incorrect*.<sup>2</sup>

**User Studies.** We conducted three IRB-approved studies to evaluate algorithm task success and querying workload: two in-lab studies, and one in-home study involving two individuals with severe mobility limitations. All studies use a Kinova Gen3 6-DoF arm with a Robotiq 2F-85 gripper, and we replicate a custom-designed feeding tool for the user studies [18]. A key methodological feature of both of our in-lab studies is the emulation of mobility constraints in participants without pre-existing mobility impairments using occupational therapy resistance bands [23].

For each method, the robot acquires 5 food items from a plate, with a maximum of 3 acquisition attempts per item to maintain a reasonable study duration. The robot may ask 4 types of queries, corresponding to each module in the bite acquisition system (Appendix C.4). The method and plate sequences are both counterbalanced across the users in each study.

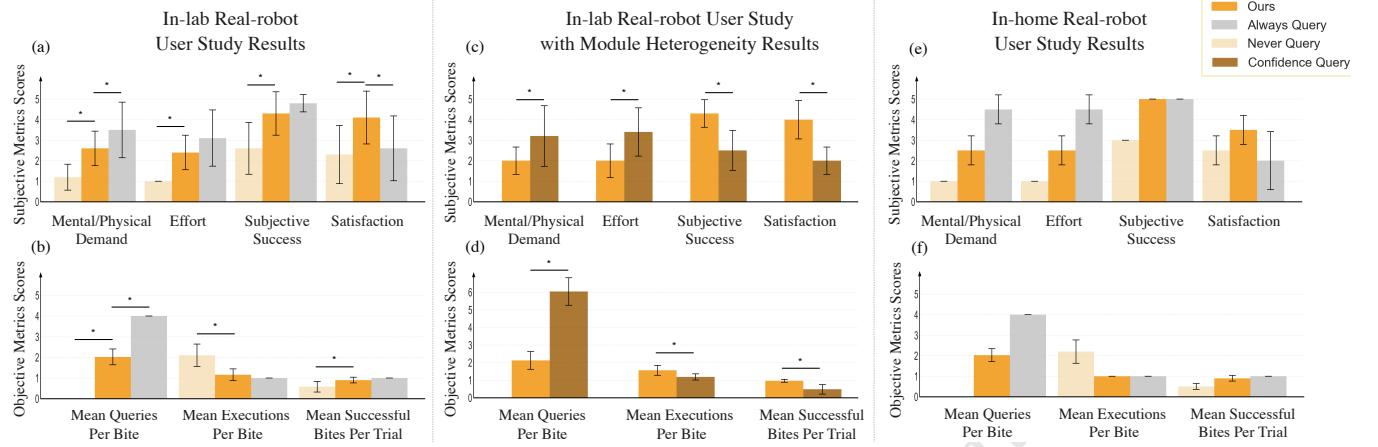
**Metrics.** We report 4 subjective metrics (Mental/Physical Demand, Effort, Subjective Success, Satisfaction) and 3 objective metrics (Mean Queries/Executions/Successful Bites Per Plate), detailed in Appendix H.

### 8.1 In-lab real-robot user study

We first conducted an in-lab study with 10 participants without mobility limitations (7 male, 3 female; ages 19–33; 40% with prior robot experience). We compared our method (`GRAPHQUERY` module selector/`QUERY-UNTIL-CONFIDENT-WORKLOAD-AWARE` querying algorithm) against two additional baselines: *Never Query* (`NEVERQUERY` module selector/`QUERY-THEN-EXECUTE` querying algorithm) and *Always Query* (`TOPOQUERY` module selector/`QUERY-FOR-ALL` querying algorithm). The study tested whether our method (1) improved task success over *Never Query* and (2) reduced querying workload relative to *Always Query*. Participants interacted with two plates: a "savory salad" (chicken, lettuce, cherry tomatoes) and a "Thanksgiving meal" (chicken, green beans, mashed potatoes).

As shown in Fig. 5(a–b), our method achieved the highest user satisfaction, significantly lower mental/physical workload than *Always Query*, and higher subjective success than *Never Query*

<sup>2</sup>While `GRAPHQUERY` and `BRUTEFORCE` were equally competitive, `BRUTEFORCE` produced false-positive queries with our binary confidence scores. We chose `QUERY-UNTIL-CONFIDENT-WORKLOAD-AWARE` over `QUERY-UNTIL-CONFIDENT` for better generalization to time-varying workload.



**Figure 5: User study metrics. In-lab real-robot study: (a) subjective scores, (b) objective scores; In-lab real-robot study with module heterogeneity: (c) subjective scores, (d) objective scores; In-home real-robot study: (e) subjective scores, (f) objective scores (\* indicates statistical significance  $p < 0.05$ ).**

(Wilcoxon signed-rank test,  $\alpha = 0.05$ ). Our method is more efficient than *Always Query*, with higher task success than *Never Query*.

## 8.2 In-lab real-robot user study with module heterogeneity

To distinguish our method from CONFIDENCE QUERY, we introduced module heterogeneity into our bite acquisition setting by incorporating a faulty food detector  $M_1$  (Appendix D), similar to our second simulation study (Sec. 7). We conducted a second in-lab study with 10 additional participants (3 male, 6 female, 1 non-binary; ages 20–30; 30% with prior robot experience) using the "savory salad" plate. We sought to determine whether the selected method improved task success rate and querying workload compared to the *Confidence Query* baseline (CONFIDENCEQUERY module selector and QUERY-UNTIL-CONFIDENT-WORKLOAD-AWARE querying algorithm).

We found that our method produced significantly higher task performance, lower workload, and higher satisfaction than *Confidence Query* (Wilcoxon test,  $\alpha = 0.05$ ; Fig. 5(c-d)), reinforcing the value of principled uncertainty and workload integration.

## 8.3 In-home real-robot user study

We conducted an in-home user study on 2 individuals with mobility limitations who require feeding assistance. One is a female, 48 years old, who has had Multiple Sclerosis for 17 years, and the other is a male, 47 years old, who has been paralyzed with a C4-C6 spinal cord injury for 27 years. Based on the in-lab study results, we compared our method to *NeverQuery* and *AlwaysQuery*. We evaluated the methods on a "mixed salad" plate with additional food item diversity (watermelons, cantaloupes, honeydew, green beans).

We find that the mental/physical demand and effort of our method are lower compared to *Always Query*, and that users rated our method as more successful compared to *Never Query* (Fig. 5(e-f)). In general, our method achieves the best user satisfaction score.

Both users expressed that they liked the system, as their satisfaction scores were higher for our method compared to both baselines. However, one user's absolute satisfaction level was reduced due to the need for physical interaction with the interface. They noted that

a more accessible option, such as voice control, would significantly improve their experience. However, they commented that because they use a tablet interface in everyday interactions, answering a few more queries on the tablet in the study did not require significantly extra effort for the study duration. However, they still rated *AlwaysQuery* as having higher effort than our method. The other user mentioned that the time required for the robot to execute actions (higher when the robot failed more frequently) was an additional source of frustration that increased overall workload.

## 9 Discussion

Our framework readily extends to other modular robot systems, which we demonstrate by considering two additional domains: a feeding architecture with visuo-haptic perceptual redundancy [39], and a large multi-robot swarm system [33]. The feeding system maps naturally to an OR-then-AND redundancy structure with low query cost ( $q_i = 0.16$ ), yielding the same recommendation as our primary setting (Fig. 4). We can model the swarm domain with  $N = 100$  to capture its scale and an OR redundancy structure, leading to CONFIDENCEQUERY for its superior scalability (Fig. 4, first column), and either QUERY-UNTIL-CONFIDENT variant. In both cases, the resulting recommendations align with domain intuition.

A key challenge in applying our framework to bite acquisition was developing well-calibrated module confidence scores. While our confidence estimates were reasonably calibrated (Sec. C.1-C.2), future work could explore conformal prediction [26, 32] or data-driven calibration methods [20]. Future work could also consider module selectors that handle more complex redundancy structures. While we studied graph-based module selectors for product-of-confidences and sum-of-uncertainties objectives, a hybrid selector could adaptively combine these structures based on which modules are redundant. Finally, longer-term interactions with end users may further differentiate human-in-the-loop algorithms, beyond our observed satisfaction trends. Future evaluations could also consider more diverse in-the-wild dishes with richer food sets, more complex geometries, and pre-manipulation skills [15, 18, 39].

## References

- [1] Christopher Agia, Rohan Sinha, Jingyun Yang, Ziang Cao, Rika Antonova, Marco Pavone, and Jeannette Bohg. 2025. Unpacking Failure Modes of Generative Policies: Runtime Monitoring of Consistency and Progress. In *Conference on Robot Learning*. PMLR, 689–723.
- [2] Faseeh Ahmad, Matthias Mayr, Sulthan Suresh-Fazeela, and Volker Krueger. 2024. Adaptable Recovery Behaviors in Robotics: A Behavior Trees and Motion Generators (BTMG) Approach for Failure Management. *arXiv preprint arXiv:2404.06129* (2024).
- [3] Rohan Banerjee, Rajat Kumar Jenamani, Sidharth Vasudev, Amal Nanavati, Sarah Dean, and Tapomayukh Bhattacharjee. 2024. To Ask or Not To Ask: Human-in-the-loop Contextual Bandits with Applications in Robot-Assisted Feeding. *arXiv preprint arXiv:2405.06908* (2024).
- [4] Suresh Bolusani, Matthieu Besançon, Ksenia Bestuzheva, Antonia Chmiela, João Dionísio, Tim Donkiewicz, Jasper van Doornmalen, Leon Eifler, Mohammed Ghannam, Ambros Gleixner, et al. 2024. The SCIP optimization suite 9.0. *arXiv preprint arXiv:2402.17702* (2024).
- [5] Maya Cakmak and Andrea L Thomaz. 2012. Designing robot learners that ask good questions. In *Proceedings of the seventh annual ACM/IEEE international conference on Human-Robot Interaction*. 17–24.
- [6] Michele Colledanchise and Petter Ögren. 2018. *Behavior trees in robotics and AI: An introduction*. CRC Press.
- [7] Cristina Cornelio and Mohammed Diab. 2024. Recover: A Neuro-Symbolic Framework for Failure Detection and Recovery. *arXiv preprint arXiv:2404.00756* (2024).
- [8] Jiafei Duan, Wilbert Pumacay, Nishanth Kumar, Yi Ru Wang, Shulin Tian, Wentao Yuan, Ranjay Krishna, Dieter Fox, Ajay Mandlekar, and Yijie Guo. 2024. Aha: A vision-language-model for detecting and reasoning over failures in robotic manipulation. *arXiv preprint arXiv:2410.00371* (2024).
- [9] Ryan Feng, Youngsun Kim, Gilwoo Lee, Ethan Gordon, Matthew Schmittle, Shivam Kumar, Tapomayukh Bhattacharjee, and Siddhartha Srinivasa. 2019. Robot-Assisted Feeding: Generalizing Skewering Strategies across Food Items on a Realistic Plate. (06 2019). doi:10.48550/arXiv.1906.02350
- [10] Tesca Fitzgerald, Pallavi Koppol, Patrick Callaghan, Russell Quinlan Jun Hei Wong, Reid Simmons, Oliver Kroemer, and Henny Admoni. 2022. INQUIRE: INteractive querying for user-aware informative REasoning. In *6th Annual Conference on Robot Learning*.
- [11] Lex Fridman, Bryan Reimer, Bruce Mehler, and William T Freeman. 2018. Cognitive load estimation in the wild. In *Proceedings of the 2018 chi conference on human factors in computing systems*. 1–9.
- [12] Émiland Garrabé, Pierre Teixeira, Mahdi Khoramshahi, and Stéphane Doncieux. 2024. Enhancing Robustness in Language-Driven Robotics: A Modular Approach to Failure Reduction. *arXiv preprint arXiv:2411.05474* (2024).
- [13] Ethan K Gordon, Xiang Meng, Tapomayukh Bhattacharjee, Matt Barnes, and Siddhartha S Srinivasa. 2020. Adaptive robot-assisted feeding: An online learning framework for acquiring previously unseen food items. In *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 9659–9666.
- [14] Ethan K Gordon, Sumegh Roychowdhury, Tapomayukh Bhattacharjee, Kevin Jamieson, and Siddhartha S Srinivasa. 2021. Leveraging post hoc context for faster learning in bandit settings with applications in robot-assisted feeding. In *2021 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 10528–10535.
- [15] Nayoung Ha, Ruolin Ye, Ziang Liu, Shubhangi Sinha, and Tapomayukh Bhattacharjee. 2024. REPeat: A Real2Sim2Real Approach for Pre-acquisition of Soft Food Items in Robot-assisted Feeding. In *2024 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 7048–7055.
- [16] Sandra G Hart and Lowell E Staveland. 1988. Development of NASA-TLX (Task Load Index): Results of empirical and theoretical research. In *Advances in psychology*. Vol. 52. Elsevier, 139–183.
- [17] Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. 2024. Gpt-4o system card. *arXiv preprint arXiv:2410.21276* (2024).
- [18] Rajat Kumar Jenamani, Priya Sundaresan, Maram Sakr, Tapomayukh Bhattacharjee, and Dorsa Sadigh. 2024. FLAIR: Feeding via Long-horizon AcQuIstion of Realistic dishes. *arXiv preprint arXiv:2407.07561* (2024).
- [19] Ulas Berk Karli, Tetsuo Kurumisawa, and Tesca Fitzgerald. [n. d.]. Ask Before You Act: Token-Level Uncertainty for Intervention in Vision-Language-Action Models. In *Second Workshop on Out-of-Distribution Generalization in Robotics at RSS 2025*.
- [20] Ulas Berk Karli, Ziyao Shangguan, and Tesca Fitzgerald. 2025. INSIGHT: INference-time Sequence Introspection for Generating Help Triggers in Vision-Language-Action Models. *arXiv preprint arXiv:2510.01389* (2025).
- [21] Ross A Knepper, Stefanie Tellex, Adrian Li, Nicholas Roy, and Daniela Rus. 2015. Recovering from failure by asking for help. *Autonomous Robots* 39 (2015), 347–362.
- [22] Shilong Liu, ZhaoYang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chunyuan Li, Jianwei Yang, Hang Su, Jun Zhu, and Lei Zhang. 2023. Grounding DINO: Marrying DINO with Grounded Pre-Training for Open-Set Object Detection. *arXiv preprint arXiv:2303.05499* (2023).
- [23] Ziang Liu, Yuanchen Ju, Yu Da, Tom Silver, Pranav N Thakkar, Jenna Li, Justin Guo, Katherine Dimitropoulou, and Tapomayukh Bhattacharjee. 2025. GRACE: Generalizing Robot-Assisted Caregiving with User Functionality Embeddings. In *2025 20th ACM/IEEE International Conference on Human-Robot Interaction (HRI)*. IEEE, 686–695.
- [24] Xiao Ma, Sumit Patidar, Iain Haughton, and Stephen James. 2024. Hierarchical Diffusion Policy for Kinematics-Aware Multi-Task Robotic Manipulation. *CVPR* (2024).
- [25] Matthias Mayr, Faseeh Ahmad, Konstantinos Chatzilygeroudis, Luigi Nardi, and Volker Krueger. 2022. Skill-based multi-objective reinforcement learning of industrial robot tasks with planning and knowledge integration. In *2022 IEEE International Conference on Robotics and Biomimetics (ROBIO)*. IEEE, 1995–2002.
- [26] James F Mullen Jr and Dinesh Manocha. 2024. Towards Robots That Know When They Need Help: Affordance-Based Uncertainty for Large Language Model Planners. *arXiv preprint arXiv:2403.13198* (2024).
- [27] Krishna Palempalli, Rohan Banerjee, Sarah Dean, and Tapomayukh Bhattacharjee. 2025. Human-in-the-loop Foundation Model Failure Recovery for Robot-Assisted Bite Acquisition. In *1st Workshop on Safely Leveraging Vision-Language Foundation Models in Robotics: Challenges and Opportunities*. <https://openreview.net/forum?id=I7UG7eC11>
- [28] Tianyang Pan, Andrew M Wells, Rahul Shome, and Lydia E Kavraki. 2022. Failure is an option: task and motion planning with failing executions. In *2022 International Conference on Robotics and Automation (ICRA)*. IEEE, 1947–1953.
- [29] Rafael Papallas and Mehmet R Dogar. 2022. To ask for help or not to ask: A predictive approach to human-in-the-loop motion planning for robot manipulation tasks. In *2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 649–656.
- [30] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. 2021. Learning Transferable Visual Models From Natural Language Supervision. *arXiv:2103.00020 [cs.CV]* <https://arxiv.org/abs/2103.00020>
- [31] Akilesh Rajavenkatnarayanan, Harish Ram Nambiappan, Maria Kyrarini, and Fillia Makedon. 2020. Towards a real-time cognitive load assessment system for industrial human-robot cooperation. In *2020 29th IEEE International Conference on Robot and Human Interactive Communication (RO-MAN)*. IEEE, 698–705.
- [32] Allen Z Ren, Anushri Dixit, Alexandra Bodrova, Sumeet Singh, Stephen Tu, Noah Brown, Peng Xu, Leila Takayama, Fei Xia, Jake Varley, et al. 2023. Robots that ask for help: Uncertainty alignment for large language model planners. *arXiv preprint arXiv:2307.01928* (2023).
- [33] Michael Rubenstein, Alejandro Cornejo, and Radhika Nagpal. 2014. Programmable self-assembly in a thousand-robot swarm. *Science* 345, 6198 (2014), 795–799.
- [34] Lucy Xiaoyang Shi, Zheyuan Hu, Tony Z Zhao, Archit Sharma, Karl Pertsch, Jianlan Luo, Sergey Levine, and Chelsea Finn. 2024. Yell at your robot: Improving on-the-fly from language corrections. *arXiv preprint arXiv:2403.12910* (2024).
- [35] Joshua Bhagat Smith, Prakash Baskaran, and Julie A Adams. 2022. Decomposed physical workload estimation for human-robot teams. In *2022 IEEE 3rd International Conference on Human-Machine Systems (ICHMS)*. IEEE, 1–6.
- [36] Priya Sundaresan, Suneel Belkhale, and Dorsa Sadigh. 2022. Learning visuo-haptic skewering strategies for robot-assisted feeding. In *6th Annual Conference on Robot Learning*.
- [37] Stefanie Tellex, Ross Knepper, Adrian Li, Daniela Rus, and Nicholas Roy. 2014. Asking for help using inverse semantics. (2014).
- [38] Shivam Vats, Michelle Zhao, Patrick Callaghan, Mingxi Jia, Maxim Likhachev, Oliver Kroemer, and George Konidaris. 2025. Optimal Interactive Learning on the Job via Facility Location Planning. *arXiv preprint arXiv:2505.00490* (2025).
- [39] Zhanxin Wu, Bo Ai, Tom Silver, and Tapomayukh Bhattacharjee. 2025. SAVOR: Skill Affordance Learning from Visuo-Haptic Perception for Robot-Assisted Bite Acquisition. *arXiv preprint arXiv:2506.02353* (2025).
- [40] Chen Xu, Tony Khuong Nguyen, Emma Dixon, Christopher Rodriguez, Patrick Miller, Robert Lee, Paarth Shah, Rares Ambrus, Haruki Nishimura, and Masha Itkina. 2025. Can we detect failures without failure data? uncertainty-aware runtime failure detection for imitation learning policies. *arXiv preprint arXiv:2503.08558* (2025).

## 1045 A Proxy objectives

1046 Justification for proxy objective 1: the second term can be a good  
 1047 approximation for  $1 - \mathbb{E}[r_{task}]$ , where the expectation is over model  
 1048 uncertainties (represented by  $c_i$ ):  
 1049

$$\begin{aligned} 1050 \mathbb{E}[r_{task}] &= P(\cap_i M_i \text{ succeeds}) \\ 1051 &= \prod_{M_i} P(M_i \text{ correct} | \text{predecessor } M_j \text{ correct}) \\ 1052 &= \prod_{M_i \notin M_q} c_i \\ 1053 &= \prod_{M_i \notin M_q} c_i \\ 1054 &= \prod_{M_i \notin M_q} c_i \\ 1055 &= \prod_{M_i \notin M_q} c_i \\ 1056 &= \prod_{M_i \notin M_q} c_i \end{aligned}$$

1057 Justification for proxy objective 2: Let  $u_i = 1 - c_i$  be the uncer-  
 1058 tainty of module  $i$ . Then we can use union bound to show that:  
 1059

$$\begin{aligned} 1060 P(\text{system fail}) &= P(\cup_i M_i \text{ fails}) \\ 1061 &\leq \sum_i P(M_i \text{ fails}) \\ 1062 &= \sum_{M_i \notin M_q} u_i + \sum_{M_i \in M_q} 0 \\ 1063 &= \sum_{M_i \notin M_q} u_i \\ 1064 &= \sum_{M_i \notin M_q} u_i \\ 1065 &= \sum_{M_i \notin M_q} u_i \\ 1066 &= \sum_{M_i \notin M_q} u_i \\ 1067 &= \sum_{M_i \notin M_q} u_i \\ 1068 &= \sum_{M_i \notin M_q} u_i \end{aligned}$$

1069 Thus, the proxy objective 2 task component is an upper bound  
 1070 in  $P(\text{system fail}) = 1 - \mathbb{E}[r_{task}]$ , meaning that proxy objective 2 is  
 1071 a loose upper bound on the original objective, so minimizing proxy  
 1072 objective 2, could also minimize the original objective.  
 1073

## B Algorithm performance analysis.

### B.1 GraphQuery

1074 In this setting, we assign edge costs  $C(e) = 1 - c_i$  for the autonomous  
 1075 edges. The graph algorithm thus queries for  $M_f$  if  $1 - c_i > \epsilon q$ ; or  
 1076

$$1077 q < (1 - c_i)/\epsilon \quad (5)$$

1078 We will never query for any of the modules in success because  
 1079  $c_i = 1$  for these modules (as  $1 - c_i$  will always be 0, and we assume  
 1080 that  $q$  is nonnegative, so condition 5 can never be met), so we don't  
 1081 have to worry about querying for modules that are upstream of  
 1082  $M_f$ .  
 1083

1084 In our scenario, we will query for  $M_f$  since  $q = 0.1$ ; and  $(1 - c_i)/\epsilon = (1 - 0.1) = 0.9$ ; thus condition 5 is met. Thus GraphQuery  
 1085 will query correctly at the first timestep.  
 1086

1087 For fixed  $q = 0.1$ , the critical value of  $\epsilon$  where we will cease to  
 1088 query is  $\epsilon_{crit} = (1 - c_i)/(q) = (1 - 0.1)/(0.1) = 9$ .  
 1089

1090 For fixed  $\epsilon = 1$ , the critical value of  $q$  above which we should  
 1091 cease to query is  $q_{crit} = (1 - c_i)/(epsilon) = 0.9$ .  
 1092

### B.2 MIP

1093 Recall that we approximate the expected task reward  $\mathbb{E}[r_{task}]$  as  
 1094 follows:  
 1095

$$\mathbb{E}[r_{task}] = \prod_{M_i} P(M_i \text{ correct}) = \prod_{M_i} c_i$$

1096 If we decide to query for a particular module  $i$ , we assume that  
 1097  $c_i = 1$  since we're using the expert feedback.  
 1098

- 1099 • Hypothetical  $J(\psi_{ms}, \psi_q)$  of not querying for any module:  
 $1 - (1^{N-1} \cdot 0.1 \cdot 1) = 0.9$
- 1100 • Hypothetical  $J(\psi_{ms}, \psi_q)$  of querying for any of the non-  
 $1101 \text{failure module: } q + 1 - (1^{N-2} \cdot 0.1 \cdot 1) = q + 0.9 = q + 0.9$
- 1102 • Hypothetical cost of querying for the failure node:  $q + 1 -$   
 $1103 (1^{N-1} \cdot 1) = q$ , so we should choose to query for the failure  
 $1104 \text{module as long as } q < 0.9$ .

1105 Thus in our scenario, MIP will always choose to query for the  
 1106 node in failure at the first timestep.  
 1107

## C Bite Acquisition Architecture

1108 Fig. 2 (left) shows the modular bite acquisition architecture that  
 1109 we use in our work, which is an adaptation of a state-of-the-art  
 1110 architecture [18], including novel foundation model components.  
 1111

1112 This autonomous system consists of four submodules, each with  
 1113 an associated confidence estimate:  
 1114

- 1115 •  $M_1$ : GPT-4o [17] food type detector that processes a whole  
 $1116 \text{plate RGB image } z_{\text{RGB}} \in \mathbb{R}^{H \times W \times 3}$  and identifies a candidate  
 $1117 \text{set of food labels } \mathcal{L} = \{l_1, l_2, \dots, l_K\}$ , where  $K$  is the number  
 $1118 \text{of unique food categories detected (e.g., cherry tomatoes,$   
 $1119 \text{pineapple, etc.). From this set, } M_1 \text{ selects the single label$   
 $1120 \text{with the highest confidence score as its output: }$

$$1121 M_1(z_{\text{RGB}}) \rightarrow l^*$$

- 1122 •  $M_2$ : GroundingDINO [22] bounding box detector that takes  
 $1123 \text{as input the selected food type label } l^* \text{ from } M_1 \text{ together with}$   
 $1124 \text{the plate RGB image } z_{\text{RGB}}. It outputs one or more bounding$   
 $1125 \text{boxes } \mathcal{B}(l^*) = \{b_1, b_2, \dots, b_J\}$ , where each  $b_j$  corresponds  
 $1126 \text{to a detected instance of the food type } l^* \text{ in the image. Thus, }$

$$1127 M_2(z_{\text{RGB}}, l^*) \rightarrow \mathcal{B}(l^*).$$

- 1128 •  $M_3$ : GPT-4o [17] skill selector that, given a detected food  
 $1129 \text{label } l^* \text{ and its corresponding bounding box } b_i$ , predicts  
 $1130 \text{the optimal skill } a_i^h \in \mathcal{A}_h$ , where  $\mathcal{A}_h$  is a set of skills (e.g.  
 $1131 \text{Skewering, Scooping, Twirling}).$

$$1132 M_3(l^*, b_i) \rightarrow a_i^h$$

- 1133 •  $M_4$ : A VLA model RT-1 that refines the skill action  $a_i^h$  into  
 $1134 \text{precise skill parameters } a_i^l \in \mathbb{R}^4$ . It takes in the skill action  
 $1135 a_i^h$  and the cropped image from the corresponding bounding  
 $1136 \text{box } b_i$ . We adapt RT-1 for modular control by fine-tuning  
 $1137 \text{it with a new regression projection head on an aggregated}$   
 $1138 \text{dataset consisting of SPANet samples [9], along with } \sim 1,000$   
 $1139 \text{additional labeled images that we collected. This design en-}$   
 $1140 \text{ables us to attribute RT-1's uncertainty purely to its skill-}$   
 $1141 \text{parameter prediction, while } M_1, M_2, \text{ and } M_3 \text{ handle food}$   
 $1142 \text{classification, location estimation, and skill selection.}$

1143 **Skill-specific parameterization.** We represent all low-  
 1144 level skills using a unified 2D action vector  $a_i^l = (x_1, y_1, x_2, y_2)$ ,  
 1145 with its semantics based on the selected skill  $a_i^h$ . For *skew-  
 1146 ering* and *twirling*, the interaction point is defined as the  
 1147 midpoint between the two predicted points,  
 1148

$$(x_c, y_c) = \left( \frac{x_1 + x_2}{2}, \frac{y_1 + y_2}{2} \right),$$

1149 1150 1151 1152 1153 1154 1155 1156 1157 1158 1159 1160

1161 and the fork tine direction  $\theta$  computed as follows:

$$1162 \quad \theta = \arctan(y_2 - y_1, x_2 - x_1).$$

1163  
1164 For *scooping*,  $(x_1, y_1)$  denotes the start of the scooping motion and  $(x_2, y_2)$  denotes the end of the scooping trajectory. This parameterization follows prior work on vision-conditioned manipulation primitives (e.g., FLAIR [18]), while enabling a unified continuous action space across heterogeneous skills.

$$1171 \quad M_4(a_i^h, b_i) \rightarrow a_i^l$$

## 1172 C.1 Calibration Datasets and Confidence 1173 Intervals

1174 To support uncertainty-aware decision-making in our bite-acquisition  
1175 pipeline, we adapt the work mentioned in [27]. We construct calibration  
1176 datasets and associated confidence intervals for all modules:  
1177  $M_1, M_2, M_3, M_4$ . Each calibration dataset records per-instance model  
1178 outputs, their confidence scores, and the corresponding ground-truth labels.

1179 **Calibration for  $M_1$ .** We use dishes and all the available plate  
1180 images from [18]. For each food type in an image,  $M_1$  generates  
1181 a ranked list of candidate tokens (labels) with log probabilities.  
1182 These are exponentiated, scaled to percentages, and rounded to two  
1183 decimal places to yield confidence scores, and the top token (token  
1184 with the highest confidence score) is selected as the predicted label  
1185 for that food type.

- 1186 • If the top token matches the ground-truth label for the food  
1187 type, we add an entry containing the top-5 tokens, their con-  
1188 fidence scores, and the ground-truth label to  $M_1$ 's calibration  
1189 dataset.
- 1190 • If the ground-truth label is not among the predictions, we  
1191 substitute a label: either an unmatched food type prediction  
1192 (if available) or, as a fallback, a matched one from the same  
1193 image. The substituted token and its associated scores are  
1194 stored alongside the ground-truth label.

1195 This substitution procedure ensures that every target food type  
1196 contributes an entry to the calibration dataset, which is neces-  
1197 sary for computing reasonable reference intervals. This yields a  
1198 confidence-based calibration dataset for  $M_1$  with 96 samples.

1199 **Calibration for  $M_2$ .** For the same set of images used in  $M_1$ , we  
1200 iterate over all bounding boxes generated by  $M_2$ . Each bounding  
1201 box is manually assessed to determine whether it is *successful* (accu-  
1202 rately capturing a target food type) or *unsuccessful* (e.g., enclosing  
1203 the wrong food type, empty regions, or background noise). For each  
1204 bounding box, we record the confidence score produced by Ground-  
1205 ingDINO. If the bounding box is deemed successful, its confidence  
1206 score is added to the success list; if unsuccessful, its confidence  
1207 score is added to the failure list. In this setup, the success scores  
1208 play the role of the “top-choice” confidence values, while the failure  
1209 scores serve as the counterpart to the “second-choice” values used  
1210 in  $M_1$  and  $M_3$ . This yields a calibration dataset of 76 samples for  
1211  $M_2$  consisting of success (61 samples) and failure (15 samples) con-  
1212 fidence distributions that are used to construct intervals analogous  
1213 to those for the other modules.

1215 **Calibration for  $M_3$ .** The identified (or substituted) label from  
1216  $M_1$ , together with the corresponding bounding box image from  $M_2$   
1217 of that food type, is passed to  $M_3$ . To ensure reliable inputs, we  
1218 restrict this step to bounding boxes deemed valid—those accurately  
1219 capturing the food type without just empty regions and such. This  
1220 module outputs candidate skill tokens and their log probabilities,  
1221 which are processed into percentage confidence scores using the  
1222 same procedure as  $M_1$ . Similar to the procedure in  $M_1$ , the top token  
1223 is taken as the predicted skill. For each food type and bounding  
1224 box image pair, we store the top 3 tokens, their confidence scores,  
1225 and the ground-truth skill label in the calibration dataset for  $M_3$ .  
1226 This yields a confidence-based calibration dataset for  $M_3$ , with 66  
1227 samples, analogous to that of  $M_1$ .

1228 **Calibration for  $M_4$ .** Since RT-1 is deterministic, we introduce  
1229 stochasticity into its predictions using Monte Carlo Dropout. Specifi-  
1230 cally, we enable dropout during inference and run the model in  
1231 batches of 16 forward passes. This produces a distribution over  
1232 the  $M_4$  output space, from which we compute the variance. The  
1233 variance then serves as our measure of the model’s confidence in its  
1234 prediction. We used the same bounding boxes as we used in  $M_3$ , but  
1235 with the human annotated ground truth label for skills. This yields  
1236 a confidence-based calibration dataset for  $M_4$  with 66 samples.

1237 **Confidence Intervals.** From each calibration dataset, we com-  
1238 pute the mean ( $\mu^*$ ) and standard deviation ( $\sigma^*$ ) of the top tokens’  
1239 confidence scores. We do the same for the second top tokens, and get  
1240  $\mu^+$  and  $\sigma^+$ . Using these values, we define the following confidence  
1241 intervals:

$$I_{\text{TopConf}} = [\mu^* - \sigma^*, \mu^* + \sigma^*], \quad I_{\text{SecondTopConf}} = [\mu^+ - \sigma^+, \mu^+ + \sigma^+].$$

1242 These intervals capture the expected distribution of confidence  
1243 scores for top and second-ranked predictions across all calibration  
1244 instances.

## 1245 C.2 Calibrated Confidence Scores

1246 To obtain calibrated confidence scores for all the modules, we apply  
1247 a rule that maps raw confidence values into a stable binary value.  
1248 For a given confidence score  $c^*$  (corresponding to the top prediction  
1249 for an instance), we check whether  $c^*$  falls outside the calibrated  
1250 top-token interval or within the second-token interval:

$$1251 \quad \text{Conf}_{\text{cal}}(c^*) = \begin{cases} 0, & \text{if } (c^* \notin I_{\text{TopConf}}) \vee (c^* \in I_{\text{SecondTopConf}}), \\ 1252 & \\ 1, & \text{otherwise.} \end{cases}$$

1253 Here, the function  $\text{Conf}_{\text{cal}}(c^*)$  produces a binary calibrated con-  
1254 fidence score: 0 when the prediction is considered low-confidence  
1255 (uncertain), and 1 when the prediction is considered high-confidence  
1256 (confident). These calibrated scores are later consumed by other  
1257 components of the pipeline that determine whether or not to query  
1258 the human for that particular module.

1259 Grounding the decision in confidence intervals estimated from  
1260 calibration data allows the system to identify both underconfident  
1261 correct predictions and overconfident incorrect predictions, pro-  
1262 ducing a more stable confidence signal than raw probabilities alone.  
1263 By introducing calibrated confidence scores, we ensure that the  
1264 system reasons over interpretable, statistically grounded intervals  
1265 rather than noisy probability values, providing more reliable and  
1266 consistent confidence estimates across all the modules.

1267

1268

1269

1270

1271

1272

1273

1274

1275

1276

1277

1278

1279

1280

1281

1282

1283

1284

1285

1286

1287

1288

1289

1290

1291

1292

1293

1294

1295

1296

1297

1298

1299

1300

1301

1302

1303

1304

1305

1306

1307

1308

1309

1310

1311

1312

1313

1314

1315

1316

1317

1318

1319

1320

1321

1322

1323

1324

1325

1326

1327

1328

1329

1330

1331

1332

1333

1334

1335

1336

1337

1338

1339

1340

### 1277 C.3 Retrieval-Augmented Feedback for Error 1278 Recovery

1280 To enable the system to recover from past errors and adapt over  
1281 time, we integrate a retrieval-augmented generation (RAG) com-  
1282 ponent into both the food type identification module ( $M_1$ ) and the  
1283 skill selection module ( $M_3$ ). The goal of this component is to incor-  
1284 porate human-provided feedback into a persistent store, allowing  
1285 the pipeline to retrieve and reuse corrections when similar inputs  
1286 are encountered in the future. This mechanism provides a means  
1287 of continual learning from mistakes and reduces repeated queries  
1288 to the human user.

1289 **Embedded Feedback Store.** We implement an EmbeddedFeed-  
1290 backStore that records feedback entries consisting of the plate image  
1291 (for  $m_1$ ) or bounding box image and food label (for  $m_3$ ), together  
1292 with the corrected output provided by the human. For  $m_1$ , embed-  
1293 dings are computed directly from the plate image using the CLIP  
1294 vision encoder [30], resulting in a purely visual representation. For  
1295  $m_3$ , embeddings are computed jointly from the bounding box im-  
1296 age of the target food item and its textual label using the CLIP  
1297 vision-language model [30], capturing both visual and semantic  
1298 context for skill selection. All embeddings are normalized and stored  
1299 persistently as vectors alongside the corresponding ground-truth  
1300 correction, which in our setup refers to the a correct food type label  
1301 for  $M_1$  and the correct skill for  $M_3$ . So whenever the system queries  
1302 the human for feedback for a particular module, the corresponding  
1303 input (image or bounding box image with label) together with its  
1304 corrected output is added to the feedback store for future retrieval  
1305 and reuse.

1306 **Retrieval.** At inference time, when  $M_1$  or  $M_3$  produces a pre-  
1307 diction, the feedback store computes the CLIP embedding for the  
1308 current input and retrieves the most similar past entry using  
1309 cosine similarity. If the best match exceeds a similarity threshold,  
1310 the stored correction is reused. This correction is converted into  
1311 a probability distribution: the retrieved food type label (or skill) is  
1312 assigned a probability proportional to the similarity score, and the  
1313 remaining probability mass is evenly distributed across the other  
1314 candidate tokens. If no sufficiently similar correction is found, the  
1315 pipeline defaults to the module's prediction.

1316 This RAG-based mechanism allows the pipeline to learn incre-  
1317 mentally from its mistakes. By leveraging similarity search, the  
1318 system avoids repeating prior errors on similar inputs and reduces  
1319 unnecessary human queries. This provides a lightweight but effec-  
1320 tive form of continual adaptation.

### 1321 C.4 Query questions and feedback in user study

1322 Four types of questions can be asked when executing the *Always*  
1323 *Query* method and our selected method in the user study:

- 1325 (1) The robot will ask the user for help by asking the following  
1326 question: "Could you tell me a food item that is on the plate?"  
1327 Users can respond by providing any valid food item that is  
1328 present on the plate.
- 1329 (2) The robot will ask the user for help by displaying the plate  
1330 image on the tablet. Users can respond by tapping 2 opposite  
1331 corners of the box on the screen to create a bounding box.
- 1332 (3) The robot will ask the user for help by asking you the fol-  
1333 lowing question: "For the food item on the plate, what skill

1335 should I use?". Users can respond by telling the robot the  
1336 skill that it should use.

- 1337 (4) The robot will ask the user for help by displaying the bound-  
1338 ing box image on the tablet. Users can respond by tapping  
1339 2 points in the following manner: for skewering, tap two  
1340 points that define the longer edge of the food item; for scoop-  
1341 ing, tap the start and the end of the scoop; for twirling, tap  
1342 two points around where the user thinks the fork should  
1343 twirl the food item.

## D Bite Acquisition Architecture with Module Heterogeneity

This architecture is similar to that described in Sec. C, with the  
following modifications:

- 1345 • M<sub>1</sub>: The food detector is faulty, producing a random food la-  
1346 bel with probability 0.9 (which lies outside of the set of food  
1347 items encountered in the "savory salad" and "Thanksgiving  
1348 meal" plates). In this scenario, the food detector assigns a  
1349 confidence of  $c_1 = 0.65$ . In the event of a success,  $M_1$  pro-  
1350 duces the label with the highest confidence score  $\ell^*$  as before,  
1351 but with a confidence score  $c_1 = \frac{c^*}{c_{max,M_1}}$ , where  $c^*$  is the  
1352 raw confidence score (Sec. C.2) and  $c_{max}$  is the maximum  
1353 raw confidence score observed in the  $M_1$  calibration set.  
1354
- 1355 • M<sub>2</sub>: Similar to Sec. C, but with a confidence score  $c_2 = \frac{c^*}{c_{max,M_2}}$ ,  
1356 where  $c^*$  is the raw confidence score (Sec. C.2) and  $c_{max,M_2}$   
1357 is the maximum raw confidence score observed in the  $M_2$   
1358 calibration set.  
1359
- 1360 • M<sub>3</sub>: Similar to Sec. C, but with a confidence score  $c_3 = \frac{c^*}{c_{max,M_3}}$ ,  
1361 where  $c^*$  is the raw confidence score (Sec. C.2) and  $c_{max,M_3}$   
1362 is the maximum raw confidence score observed in the  $M_3$   
1363 calibration set.  
1364
- 1365 • M<sub>4</sub>: Unchanged compared to Sec. C.  
1366

## E Workload Model Conditioning

To predict workload for each of the bite acquisition modules (Appendix C), we set the query type variables for the predictive workload model as follows [3]:

- 1367 • M<sub>1</sub>:  $d_t = \text{easy}$ ,  $resp_t = \text{MCQ}$ ,  $dist_t = \text{"no distraction task"}$
- 1368 • M<sub>2</sub>:  $d_t = \text{easy}$ ,  $resp_t = \text{BB}$ ,  $dist_t = \text{"no distraction task"}$
- 1369 • M<sub>3</sub>:  $d_t = \text{easy}$ ,  $resp_t = \text{MCQ}$ ,  $dist_t = \text{"no distraction task"}$
- 1370 • M<sub>4</sub>:  $d_t = \text{hard}$ ,  $resp_t = \text{BB}$ ,  $dist_t = \text{"no distraction task"}$

## F Additional Synthetic Experiments: Systematic Ablations

### F.1 Additional information about experimental setup.

We set the GraphQuery hyperparameter  $\epsilon = 1$ , and we assume equal weight between querying and task reward ( $w = 0.5$ ).

Independent variable settings considered:

- 1383 • Number of modules  $N$ : [3, 5, 10, 15, 18, 25, 50, 75, 100]
- 1384 • Graph redundancy structure  $G_M$ : [fully redundant (all-OR),  
1385 fully dependent (all-AND), fully redundant followed by fully  
1386 dependent (OR-then-AND), fully dependent followed by fully  
1387 redundant (AND-then-OR)]

1388  
1389  
1390  
1391  
1392

- 1393 • Module confidences  $c_i$ :  $[(1.0, 0.1), (0.9, 0.2), (0.8, 0.3), (0.7,$   
 1394  $0.4)]$ , where  $(c_h, c_l)$  are the high confidence value and low  
 1395 confidence value, respectively, described in Sec. 6. We assume  
 1396 that 3 of the modules in the module graph  $G_M$  are assigned  
 1397 the low confidence value, and the rest are assigned the high  
 1398 confidence value.  
 1399 • Query costs  $q_i$ :  $[0.08, 0.16, 0.32, 0.64]$

## F.2 Number of modules.

Additional takeaways:

- 1400 • No other metrics vary as a function of number of modules,  
 1401 except in different variable settings, (e.g. closer confidence  
 1402 values) where numerical stability is a concern.  
 1403 • When confidence values differ from  $(1, 0.1)$ , selectors de-  
 1404 grade at larger graph sizes due to numerical underflow. We  
 1405 find that Binary Tree Query is the most sensitive, followed  
 1406 by Graph query, followed by Brute Force.

## F.3 Redundancy structures

### F.3.1 Module selectors.

- 1411 • **Query Cost.** We generally observe that the query cost is  
 1412 lower for the all-OR redundancy structure, compared to  
 1413 the other redundancy structures, because fewer queries are  
 1414 needed for the overall system output to be correct. Addi-  
 1415 tionally, the query cost for the OR-then-AND structure is slightly  
 1416 higher than that for the AND-then-OR redundancy structure  
 1417 (because overall system success is more likely when the final  
 1418 module is an OR module).  
 1419 • **Total Failed Attempts.** Again, we observe that the Total  
 1420 Failed Attempts is lower for the all-OR redundancy structure.  
 1421 The relative ordering of module selectors does not generally  
 1422 depend on the redundancy structure, and proceeds roughly  
 1423 as follows (in descending order): Never Query, Binary Tree  
 1424 Query, Brute Force, followed by Graph Query.  
 1425 • **Computation Time.** Again, we observe that the Compu-  
 1426 tation Time is lower for the all-OR redundancy structure.  
 1427 GraphQuery tends to have a higher Computation Time than  
 1428 the other module selectors (regardless of the redundancy  
 1429 structure) due to the computational complexity of parallel  
 1430 graph creation + parallel shortest-path searches.  
 1431 • **Total Correct.** For some experimental settings, we have  
 1432 complete success for all module selectors, regardless of de-  
 1433 pendency structure. For other settings, Never Query and  
 1434 Binary Tree Query only achieve complete success for the  
 1435 all-OR architecture, and have no success for the other redun-  
 1436 dancies.

### F.3.2 Querying algorithms.

- 1437 • **Query Cost.** Lower for all querying algorithms in all-OR  
 1438 redundancy structure, compared to other structures. In low  
 1439 query cost settings: (1) higher for QUERY-THEN-EXECUTE-  
 1440 WORKLOAD-AWARE compared to other methods, (2) higher  
 1441 in the hybrid structures (compared to the pure structures)  
 1442 across-the-board.  
 1443 • **Total Failed Attempts.** Lower for all querying algorithms  
 1444 in all-OR redundancy structure, compared to other structures.

Notice the following roughly consistent ordering across  
 1445 querying algorithms, regardless of graph structure (descend-  
 1446 ing order): EXECUTE-FIRST > QUERY-THEN-EXECUTE > QUERY-  
 1447 UNTIL-CONFIDENT QUERY-UNTIL-CONFIDENT-WORKLOAD-  
 1448 AWARE, which is intuitive. Notice the following ordering  
 1449 across graph structures, regardless of querying algorithm  
 1450 (descending order): OR-then-AND, all-AND, AND-then-OR,  
 1451 all-OR.

- 1452 • **Computation Time.** Lower for all querying algorithms in  
 1453 all-OR redundancy structure, compared to other structures.  
 1454 Noticed that QUERY-THEN-EXECUTE-WORKLOAD-AWARE has  
 1455 slightly higher computation time than other querying algo-  
 1456 rithms (independent of graph structure).  
 1457 • **Total Correct.** Fairly consistent across graph structures.  
 1458 QUERY-THEN-EXECUTE-WORKLOAD-AWARE has the highest  
 1459 success rate across all querying algorithms (able to recover  
 1460 in some settings where all methods fail, for AND-then-OR,  
 1461 and OR-then-AND).  
 1462 • **Total Timesteps.** Lower for all querying algorithms in all-  
 1463 OR redundancy structure, compared to other structures. Notice  
 1464 the following roughly consistent ordering across query-  
 1465 ing algorithms, regardless of graph structure (descending  
 1466 order): EXECUTE-FIRST > QUERY-THEN-EXECUTE > QUERY-  
 1467 UNTIL-CONFIDENT QUERY-UNTIL-CONFIDENT-WORKLOAD-  
 1468 AWARE, lining up with the Total Failed Attempts trends.

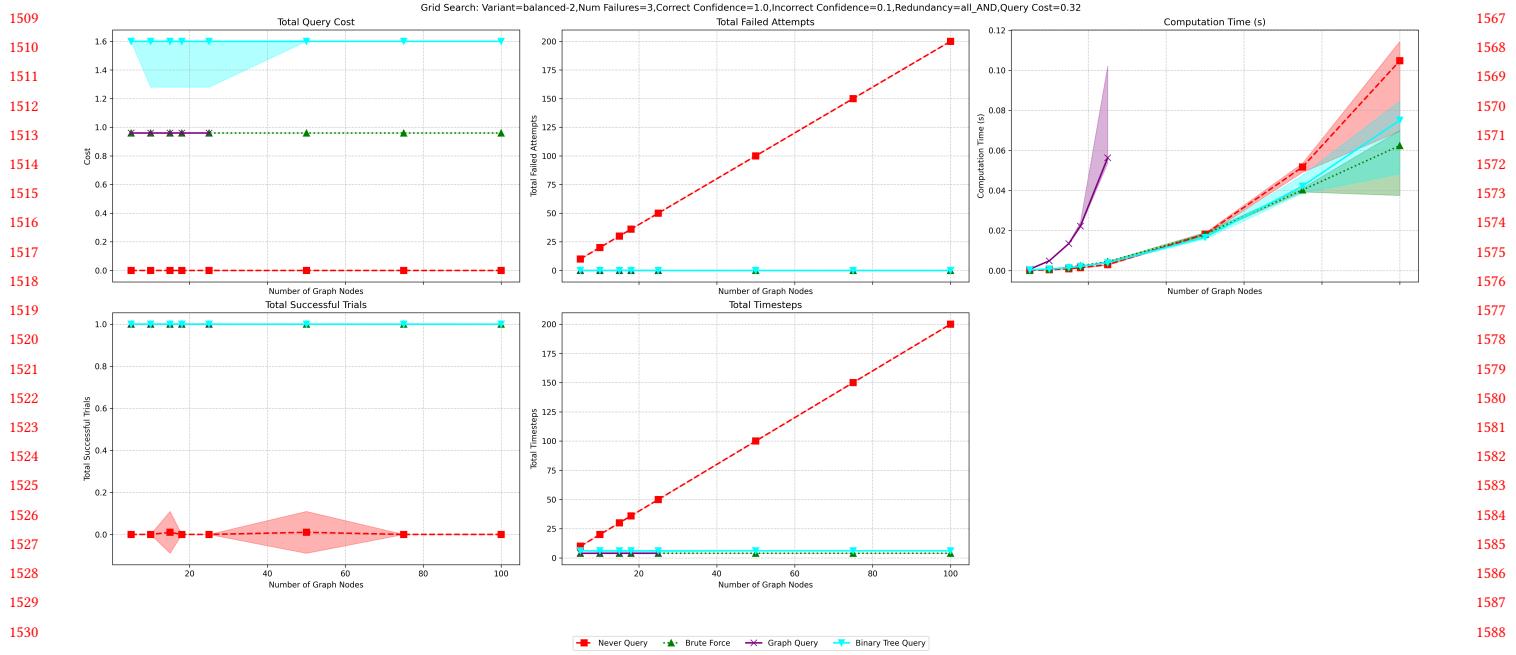
## F.4 Confidences

### F.4.1 Module selectors.

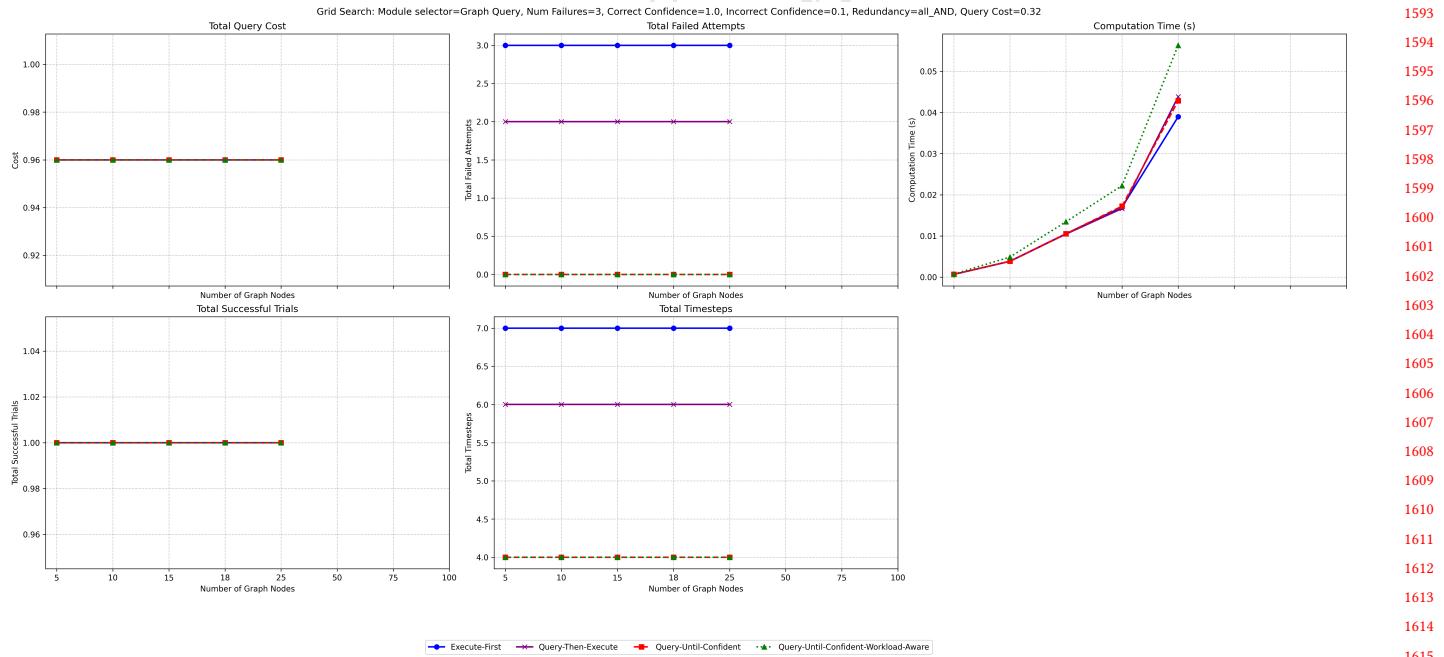
- 1469 • **Query Cost.** Generally higher when confidences are closer  
 1470 together. Regardless of confidence level, generally follows  
 1471 the following pattern in descending order: Binary Tree Query,  
 1472 Graph Query Brute Force, Never Query.  
 1473 • **Total Failed Attempts.** Generally higher when confidences  
 1474 are closer together. Regardless of confidence level, generally  
 1475 follows the following pattern in descending order: Never  
 1476 Query, Binary Tree Query, Graph Query Brute Force.  
 1477 • **Computation Time.** Generally higher when confidences  
 1478 are closer together. Graph Query has much higher Compu-  
 1479 tation Time, followed by Graph Query Brute Force, followed  
 1480 by Never Query.  
 1481 • **Total Correct.** Relatively consistent at 1.0 for Brute Force  
 1482 and Graph Query (regardless of confidence setting), while  
 1483 Binary Tree Query tends to drop from 1.0 to 0.0 with higher  
 1484 confidences. Never Query always has a value of 0.0, except  
 1485 in the all-OR setting.  
 1486 • **Total Timesteps.** Generally higher when confidences are  
 1487 closer together (for all module selectors except Never Query).  
 1488 Brute Force becomes relatively worse compared to the other  
 1489 module selector in closer confidence settings, although rank-  
 1490 ing is generally Binary Tree Query followed by Brute Force/Graph  
 1491 Query.

### F.4.2 Querying algorithms.

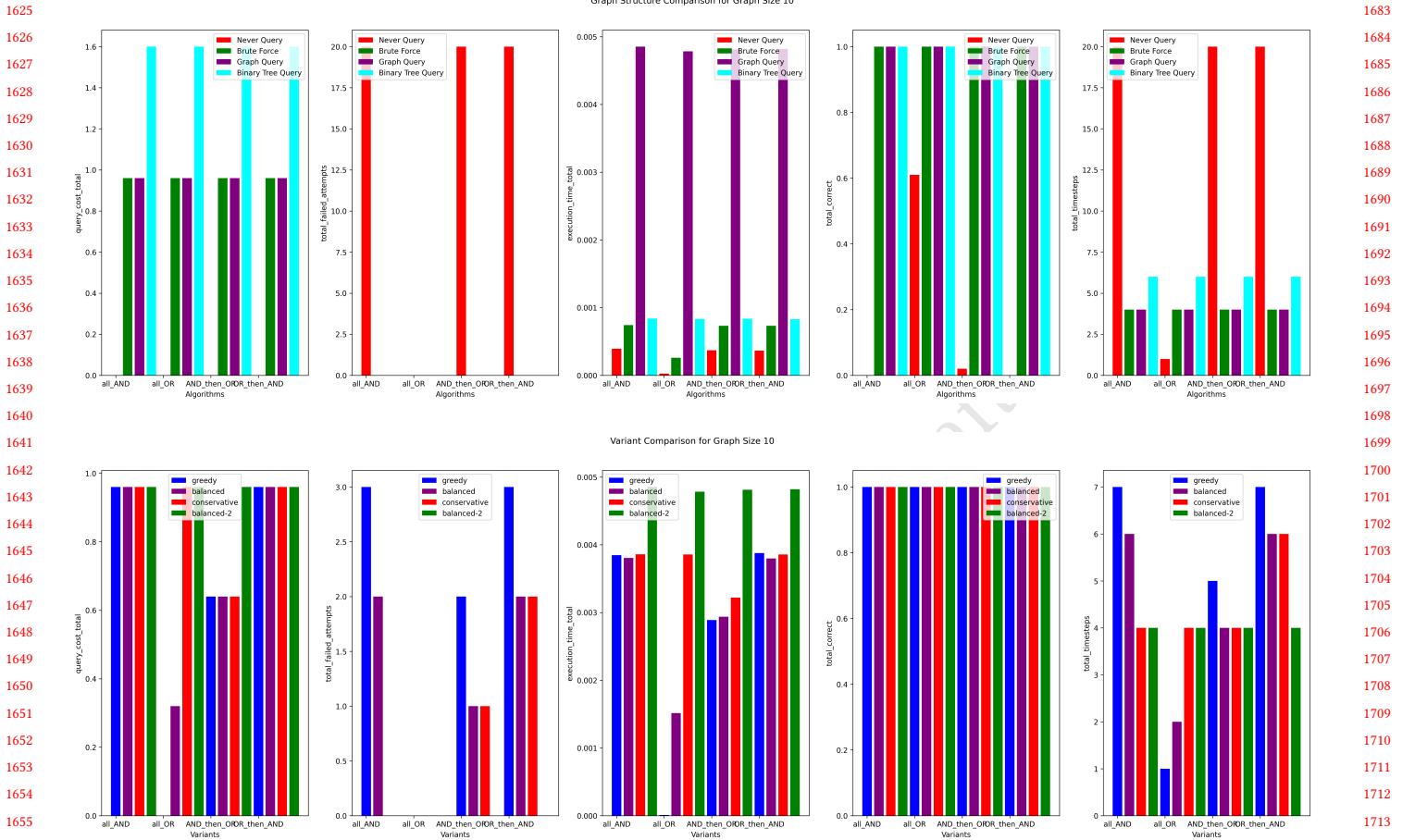
1492 1493 1494 1495 1496 1497 1498 1499 1500 1501 1502 1503 1504 1505 1506 1507 1508



**Figure 6: Varying the number of modules in the module graph  $G_M$ , for the QUERY-UNTIL-CONFIDENT-WORKLOAD-AWARE variant. Plots show median, upper quartile, and lower quartile values across 100 trials.**



**Figure 7: Varying the number of modules in the module graph  $G_M$ , for the GRAPH-QUERY querying algorithm. Plots show median, upper quartile, and lower quartile values across 100 trials.**



**Figure 8: Varying the graph redundancy structure for fixed querying algorithm (QUERY-UNTIL-CONFIDENT-WORKLOAD-AWARE; top) and fixed module selector (BRUTE-FORCE; bottom). Plots show median values across 100 trials.**

- **Query Cost.** Generally higher when confidences are closer together. Generally higher when confidences are closer together. For low query cost settings, QUERY-UNTIL-CONFIDENT-WORKLOAD-AWARE has the highest Total Query Cost compared to other querying algorithms, but for higher query cost settings it tends to have the lowest Total Query Cost.
- **Total Failed Attempts.** Generally higher when confidences are closer together. Lower for QUERY-UNTIL-CONFIDENT-WORKLOAD-AWARE and QUERY-UNTIL-CONFIDENT than EXECUTE-FIRST and QUERY-THEN-EXECUTE.
- **Computation Time.** Generally higher when confidences are closer together. Slightly higher for QUERY-UNTIL-CONFIDENT-WORKLOAD-AWARE compared to the other methods.
- **Total Correct.** Relatively consistent (1.0) for all querying algorithms, although EXECUTE-FIRST and QUERY-THEN-EXECUTE tend to decay to 0.0 for a higher number of un-confident modules.
- **Total Timesteps.** Generally higher when confidences are closer together. Regardless of confidence score, Total Timesteps

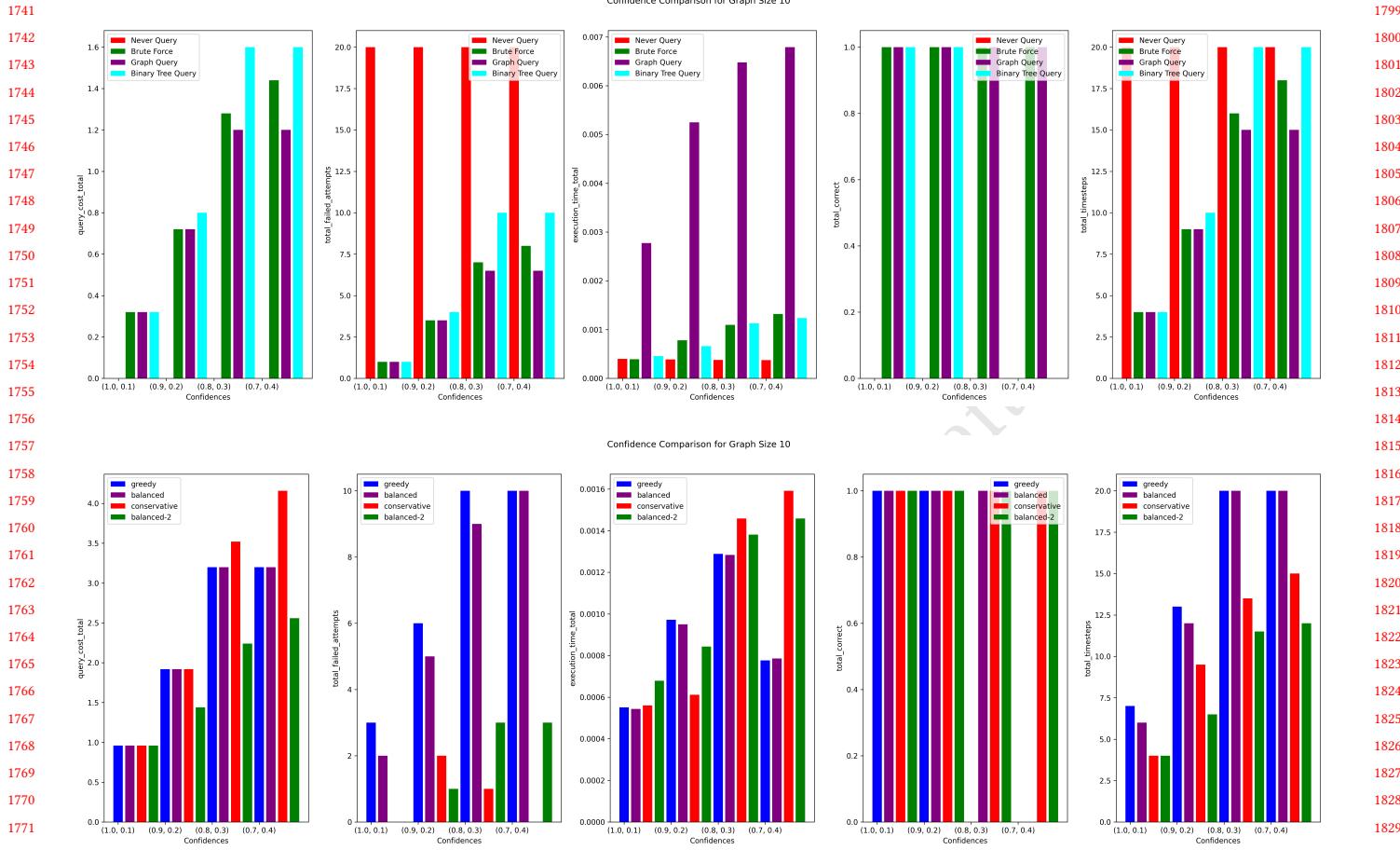
is generally lower for QUERY-UNTIL-CONFIDENT-WORKLOAD-AWARE and QUERY-UNTIL-CONFIDENT compared to EXECUTE-FIRST and QUERY-THEN-EXECUTE (except for the all-OR redundancy structure, where the trend is flipped).

*Recommendations.* For module selectors, we recommend using either Brute Force or Graph Query, as they are the most competitive across metrics regardless of the confidence level. For querying algorithms, we recommend QUERY-UNTIL-CONFIDENT-WORKLOAD-AWARE or QUERY-UNTIL-CONFIDENT (unless in the all-OR setting, in which case EXECUTE-FIRST or QUERY-THEN-EXECUTE are the best at minimizing Total Timesteps).

## F.5 Query Cost

### F.5.1 Module selectors.

- **Query Cost.** Higher total query cost across-the-board as we increase the cost of querying, which is expected. Binary Tree Query generally has the highest Total Query Cost (regardless of module query cost), while Brute Force and Graph Query are lower.



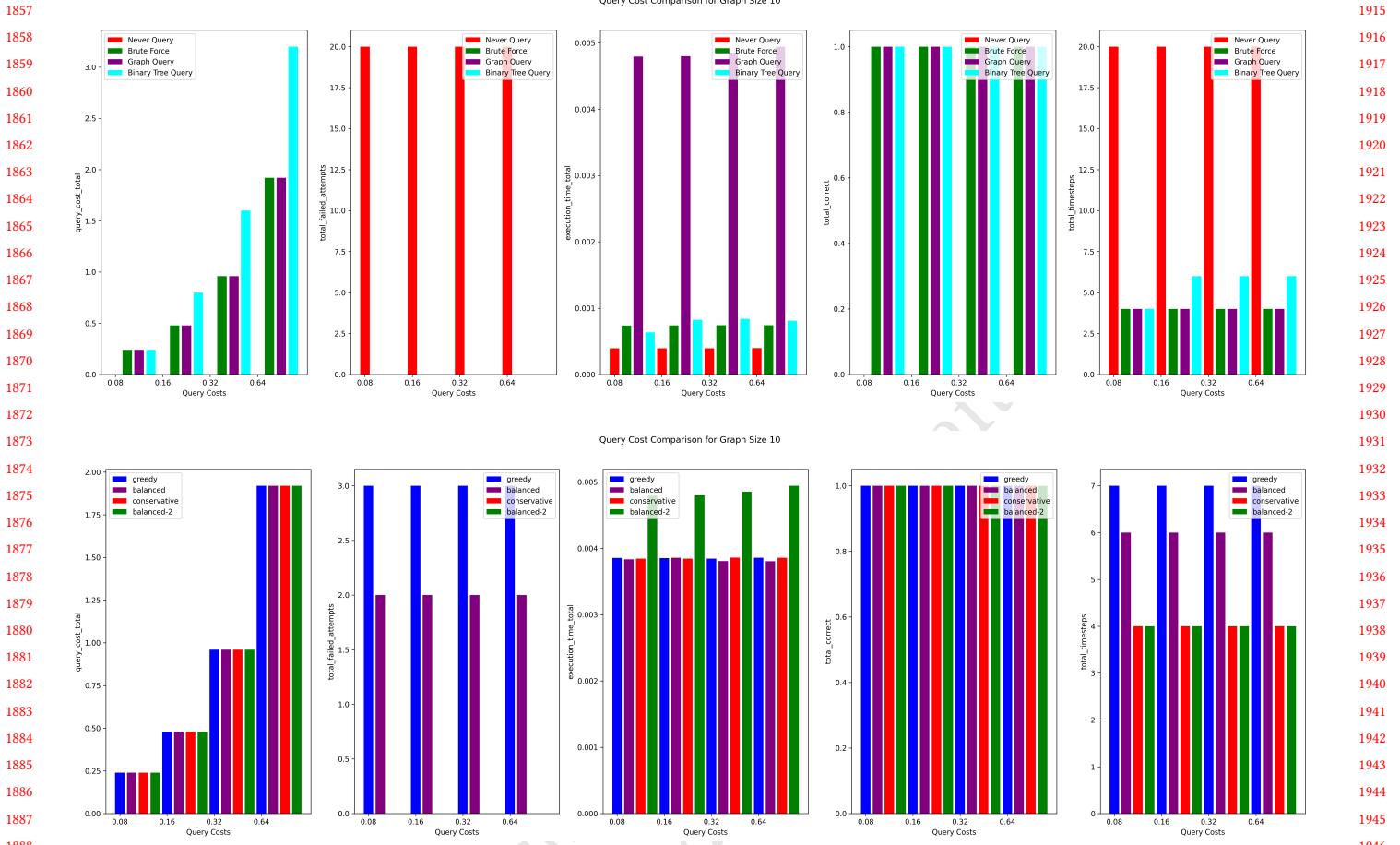
**Figure 9: Varying the confidence levels for fixed querying algorithm (QUERY-THEN-EXECUTE; top) and fixed module selector (BRUTE-FORCE; bottom). Plots show median values across 100 trials.**

- **Total Failed Attempts.** Relatively invariant to the query cost level.
- **Computation Time.** Generally highest for Graph Query, regardless of the query cost level. In ascending order for the other module selectors, generally see Never Query, then Brute Force, then Binary Tree Query.
- **Total Correct.** Generally consistent for all methods, regardless of query cost level. For higher query cost settings and high failure/low confidence settings, we do observe consistent a Total Correct value of 0.0 across all settings for Binary Tree Query and Never Query.
- **Total Timesteps.** Values are also generally invariant as a function of query cost, as all of the module selectors are deciding to query regardless of the query cost (since we are examining performance with the QUERY-THEN-EXECUTE querying algorithm).

### F.5.2 Querying algorithms.

- **Query Cost.** Total Query Cost generally goes up, but as query cost increases, we start to see a larger separation between QUERY-UNTIL-CONFIDENT-WORKLOAD-AWARE (which is lower) and the other 3 querying algorithms.
- **Total Failed Attempts.** Typical trend is EXECUTE-FIRST > QUERY-THEN-EXECUTE > QUERY-UNTIL-CONFIDENT > QUERY-UNTIL-CONFIDENT-WORKLOAD-AWARE, regardless of the query cost setting. In some settings, Total Failed Attempts increases slightly for QUERY-UNTIL-CONFIDENT-WORKLOAD-AWARE as a function of query cost.
- **Computation Time.** Generally comparable across querying algorithms, invariant to query cost setting. In some scenarios, QUERY-UNTIL-CONFIDENT-WORKLOAD-AWARE has a higher computation time than the other 3 querying algorithms.
- **Total Correct.** Generally comparable across querying algorithms (1.0), invariant to query cost setting. Sometimes observe degradation to 0.0 (in low confidence/high number of unconfident module settings) for all methods besides QUERY-UNTIL-CONFIDENT-WORKLOAD-AWARE.

1856



**Figure 10: Varying the module query costs for fixed querying algorithm (QUERY-THEN-EXECUTE; top) and fixed module selector (BRUTE-FORCE; bottom). Plots show median values across 100 trials.**

- **Total Timesteps.** Generally invariant to query cost setting. Regardless of confidence score, Total Timesteps is generally lower for QUERY-UNTIL-CONFIDENT-WORKLOAD-AWARE and QUERY-UNTIL-CONFIDENT compared to EXECUTE-FIRST and QUERY-THEN-EXECUTE (except for the all-OR redundancy structure, where the trend is flipped).

*Recommendations.* For module selectors, we recommend using either Brute Force or Graph Query, as they are the most competitive across metrics regardless of the module query cost. For querying algorithms, we recommend QUERY-UNTIL-CONFIDENT-WORKLOAD-AWARE or QUERY-UNTIL-CONFIDENT (QUERY-UNTIL-CONFIDENT-WORKLOAD-AWARE if minimizing Total Query Cost and Total Failed Attempts is most important; QUERY-UNTIL-CONFIDENT if minimizing Computation Time is the most important).

## G Additional Synthetic Experiments: Module Heterogeneity

Fig. 11 compares the GRAPHQUERY vs CONFIDENCEQUERY module selectors for additional values of the module confidences and 2026-01-02 02:45. Page 17 of 1-18.

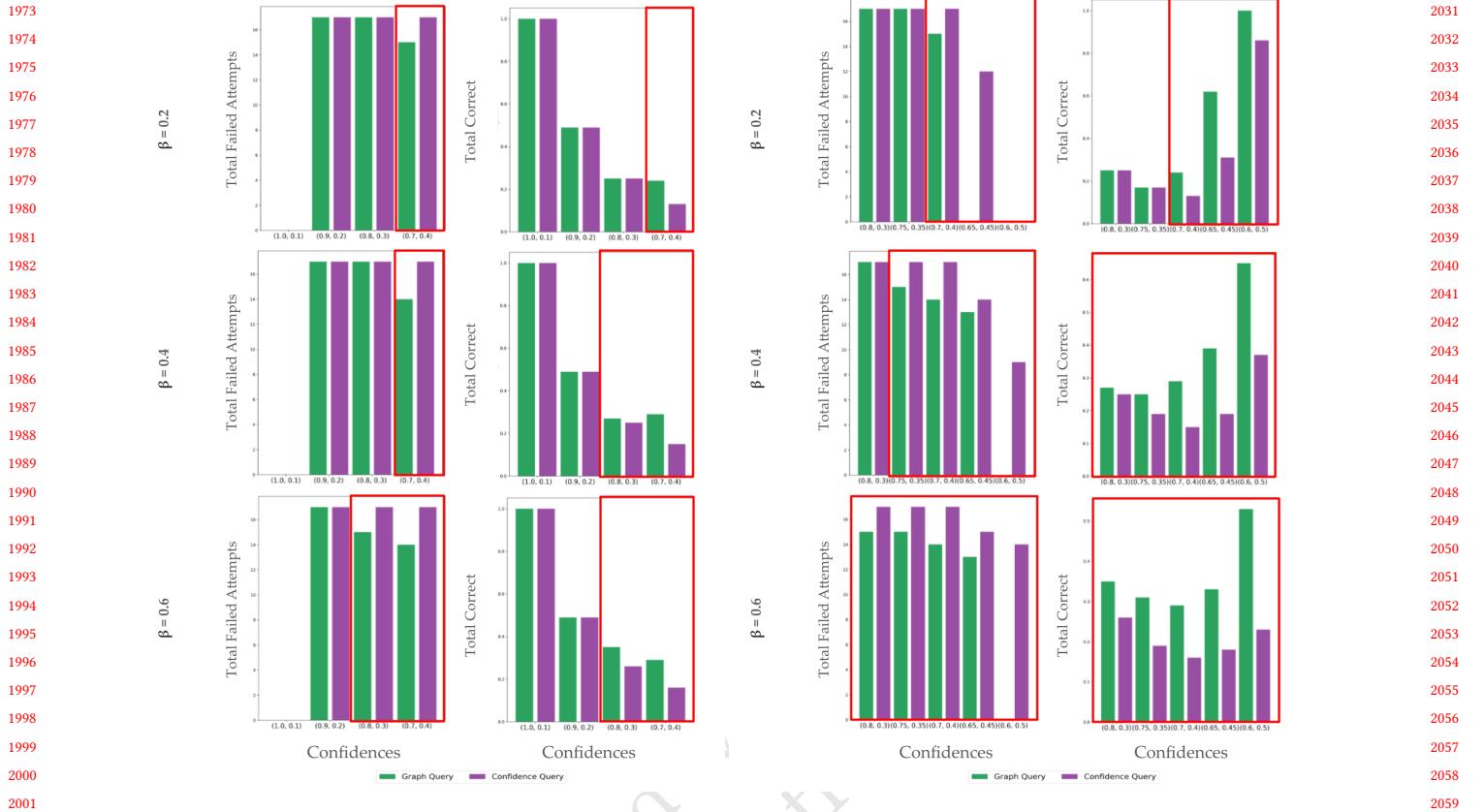
query cost variance  $\beta$  (with fixed querying algorithm QUERY-UNTIL-CONFIDENT-WORKLOAD-AWARE). We find that GRAPHQUERY performs at least as well as CONFIDENCEQUERY.

## H User Study Metrics

### H.1 Subjective metrics

We asked the participants the following questions in terms of Mental/Physical Demand, Effort, Subjective Success, and Satisfaction, all on a Likert scale from 1-5:

- (1) **Mental/Physical Demand:** For the last method, how mentally/physically demanding was it for the robot to query you?
- (2) **Effort:** For the last method, how hard did you have to work to make the robot pick up food items?
- (3) **Subjective Success:** For the last method, how successful was the robot in picking up food items?
- (4) **Satisfaction:** For the last method, how satisfied are you with how the robot balanced between trying to pickup independently when possible and asking for help when required?



**Figure 11: Comparison between GRAPHQUERY and CONFIDENCEQUERY, as a function of varying module confidences and query cost uniform noise  $\beta$  (for fixed querying algorithm: QUERY-UNTIL-CONFIDENT-WORKLOAD-AWARE). Scenarios where GRAPHQUERY outperforms CONFIDENCEQUERY are highlighted in red, both in less overlapping (left) and more overlapping (right) confidence regimes.**

## H.2 Objective metrics

We define the following 3 objective metrics:

- (1) **Mean Queries Per Plate:**  $\frac{1}{B} \sum_{b=1}^B \sum_{t=1}^{T_b} \mathbb{1}\{\text{queried at } t\}$ , where  $B$  is the number of bites per plate and  $T_b$  is the number of timesteps needed for bite  $b$ .

- (2) **Mean Executions Per Plate:**  $\frac{1}{B} \sum_{b=1}^B \sum_{t=1}^{T_b} \mathbb{1}\{\text{executed at } t\}$ , where  $B$  and  $T_b$  are defined above.
- (3) **Mean Successful Bites Per Plate:**  $\frac{1}{B} \sum_{b=1}^B \mathbb{1}\{\text{b succeeded}\}$ , where  $B$  is defined above and we define bite  $b$  to have succeeded if the robot acquired the bite after  $T_b$  timesteps.