

Large-Language-Model-Guided State Estimation for Partially Observable Task and Motion Planning

Yoonwoo Kim¹, Raghav Arora¹, Roberto Martín-Martín¹, Peter Stone^{1,3},
Ben Abbatematteo^{1†}, Yoonchang Sung^{2†}

Abstract—Robot planning in partially observable environments, where not all objects are known or visible, is a challenging problem, as it requires reasoning under uncertainty through partially observable Markov decision processes. During the execution of a computed plan, a robot may unexpectedly observe task-irrelevant objects, which are typically ignored by naive planners. In this work, we propose incorporating two types of common-sense knowledge: (1) certain objects are more likely to be found in specific locations; and (2) similar objects are likely to be co-located, while dissimilar objects are less likely to be found together. Manually engineering such knowledge is complex, so we explore leveraging the powerful common-sense reasoning capabilities of large language models (LLMs). Our planning and execution framework, CoCo-TAMP, introduces a hierarchical state estimation that uses LLM-guided information to shape the belief over task-relevant objects, enabling efficient solutions to long-horizon task and motion planning problems. In experiments, CoCo-TAMP achieves an average reduction of 62.7% in planning and execution time in simulation, and 72.6% in real-world demonstrations, compared to a baseline that does not incorporate either type of common-sense knowledge.

I. INTRODUCTION

Robots performing long-horizon manipulation tasks must reason over discrete decisions, such as which objects to interact with, and continuous motions for manipulation and navigation. Task and motion planning (TAMP) [1] provides a principled approach for such problems. However, in realistic settings with uncertainty over object poses and occlusions, plans from deterministic TAMP solvers can fail. This work addresses the partially observable task and motion planning (PO-TAMP) problem, aiming to enable robots to effectively plan to manipulate objects that may not be directly visible due to partial observability.

We introduce CoCo-TAMP, a TAMP system that uses large language models (LLMs) to provide common-sense priors and co-location cues that shape beliefs during planning and execution. While LLMs have limitations as planners [2], [3], they can serve as a rich, approximate source of knowledge [4], [5]. LLM-Modulo [5] provides a framework that leverages the strength of LLMs with external model-based verifiers in a “generate and verify” loop, allowing flexibility in problem specification while guaranteeing the completeness of the framework. Following this paradigm, CoCo-TAMP queries an LLM to form priors over rooms/surfaces via multiple-choice question answering (MCQA) and uses LLM sentence embeddings to build a similarity-based co-location

model that propagates evidence across objects. Beliefs are maintained in a hierarchical Bayesian filter over rooms, surfaces, and continuous poses, with a visibility-aware observation model that addresses misses under partial coverage and a lightweight “co-location toggler” that disables co-location when semantics suggest broad dispersion (e.g., switches). Integrated with a belief-space planner, these components yield efficient information gathering and execution. In simulation and on a real robot, CoCo-TAMP reduces cumulative planning and execution time by 62.7% and 72.6%, respectively, compared to a baseline without LLM priors or co-locations.

In summary, this paper makes the following contributions:

- We propose an interleaved planning–execution framework for PO-TAMP that leverages LLMs to provide commonsense knowledge. By providing informative priors and guiding belief updates, CoCo-TAMP enables more accurate object state estimation and makes belief-space planning practical under partial observability.
- We demonstrate the effectiveness of our approach through large-scale household simulations and real-world robot experiments. Across diverse environments, our framework substantially reduces planning and execution time compared to baselines, and maintains robustness even in adversarial settings where common-sense priors are misleading.

II. PROBLEM DESCRIPTION

We assume a known semantic layout (rooms and surfaces) obtained from semantic SLAM system [6], but only partial information about objects is available. The environment contains O manipulable objects, R rooms, and S surfaces. For each task-relevant object $k \in \{1, \dots, K\}$, let $x_{r,t}^k \in \mathcal{R}$, $x_{s,t}^k \in \mathcal{S}$, and $x_{p,t}^k \in SE(3)$ denote its room, surface, and pose at time t , with beliefs $\text{bel}(x_{r,t}^k)$, $\text{bel}(x_{s,t}^k)$ be categorical and $\text{bel}(x_{p,t}^k)$ be continuous. We are given initial beliefs $\text{bel}(x_{r,0}^k)$, $\text{bel}(x_{s,0}^k)$, and $\text{bel}(x_{p,0}^k)$ over $K < O$ task-relevant objects.

Building on SS-Replan [7], we model PO-TAMP as a hybrid discrete–continuous belief-space stochastic shortest-path problem [8] where actions have strictly positive cost and the minimum-cost plan maximizes success likelihood. CoCo-TAMP takes a TAMP specification $(\mathcal{O}, \mathcal{P}, \mathcal{I}, \mathcal{G}, \mathcal{A})$, where \mathcal{O} is the set of manipulable objects, \mathcal{P} the predicates, \mathcal{I} the initial literals, \mathcal{G} the goal literals, and \mathcal{A} a finite set of parameterized actions with preconditions and effects. Throughout planning and execution, CoCo-TAMP maintain

¹The University of Texas at Austin, ²Nanyang Technological University, ³Sony AI, [†]Equal advising
Correspondence: yoonchang.sung@ntu.edu.sg

beliefs over $(x_{r,t}^k, x_{s,t}^k, x_{p,t}^k)$ for each task-relevant object k and triggers replanning whenever execution failures occur.

PDDLStream [9] was used as the underlying TAMP planner to couple symbolic actions with streams that sample and certify continuous variables (*e.g.*, poses and trajectories). In addition to standard pick, place, and move actions, we include an observation action `detect`. The `detect` action is parameterized by the object `?o`, surface `?s`, room `?r`, sampled object pose from $\text{bel}(x_{p,0}^k)$ `?pb`, base configuration `?bq`, head configuration `?hq`, and head trajectory `?ht` of the robot. Following the self-loop MDP heuristic in SS-Replan, the cost of taking the `detect` action (`DetectCost`) is set inversely to the belief over the object’s state, which serves as a proxy for detection likelihood (*i.e.*, cheaper when belief mass and visibility predict success). Concretely, this cost is computed as $\frac{c_a}{\text{bel}(x_{r,t}^k, x_{s,t}^k, x_{p,t}^k)}$ where $c_a = 1$ and the resulting cost steers the planner toward informative views.

We are particularly interested in improving the efficiency of solving PO-TAMP problems by achieving accurate beliefs guided by LLM-driven initializations and object co-location models. To this end, we evaluate COCO-TAMP using two metrics: the cumulative planning and execution time until successful task completion, and the number of replanning iterations required to solve the problem. The latter metric is motivated by the fact that the former can vary significantly depending on implementation details, system specifications, and the task itself. As such, the number of execution failures serves as a useful proxy metric for evaluating efficiency.

III. COCO-TAMP

The overview of COCO-TAMP is presented in Figure 1, which illustrates an integrated planning and execution framework. LLM-guided initial belief generation is discussed in Section III-A. Section III-B describes how the beliefs are updated during execution using a hierarchical Bayes filter with LLM-informed semantic similarity between objects.

A. Initial belief generation with LLMs

LLMs trained on internet-scale data possess common-sense knowledge about typical object placements. Leveraging this capability, COCO-TAMP utilizes LLMs to generate prior beliefs over rooms and surfaces, denoted by $\text{bel}(x_{r,0}^k)$ and $\text{bel}(x_{s,0}^k)$.

Prior work [10], [11], [12] has employed MCQA for uncertainty quantification in LLMs. We adopt their approach by formulating the most likely room and surface selection problem as an MCQA task. Possible rooms and surfaces where the object could be located are paired with single-character options (*i.e.*, A, B, C, ...). The LLM is asked to choose the character corresponding to the room and surface where the object is most likely to be located. An example input prompt given to the LLM to estimate $\text{bel}(x_{r,0}^k)$ is presented below, and the corresponding prompt to estimate $\text{bel}(x_{s,0}^k)$ omitted, as it follows an analogous structure.

Predict the location of a toaster.

- (A) Kitchen
- (B) Bathroom

(C) Livingroom

(D) Garage

Return the letter that represents the location:

The probabilities over potential answers can be derived from the LLM’s next-token prediction mechanism by evaluating the log probabilities of each possible answer’s tokens when appended to the prompt. Consequently, given a one-shot input prompt similar to the one above, we obtain logit scores corresponding to the option choices (*i.e.*, A, B, C, ...). These logit scores are then passed through a softmax function, which converts them into valid probabilities.

B. Object state estimation

To obtain the probability of successful detection, we jointly estimate the semantic location (*i.e.*, surface, room) and the pose of the objects while accounting for the *visibility* of the location. Formally, we aim to compute the posterior probability of the k -th object’s semantic location and pose, denoted as $\text{bel}(x_{r,t}^k, x_{s,t}^k, x_{p,t}^k)$, using the recursive Bayesian filtering framework.

Objects are on surfaces, and surfaces are in rooms. Consequently, the belief distribution $\text{bel}(x_{r,t}^k, x_{s,t}^k, x_{p,t}^k)$ is factored into three conditional terms:

- $\text{bel}(x_{p,t}^k) = P(x_{p,t}^k | x_{s,t}^k, x_{r,t}^k, z_{p,t}^k)$
- $\text{bel}(x_{s,t}^k) = P(x_{s,t}^k | x_{r,t}^k, \mathcal{Z}_{s,t}, v_{s,t})$
- $\text{bel}(x_{r,t}^k) = P(x_{r,t}^k | \mathcal{Z}_{r,t}, v_{r,t})$

The observation models include both continuous and categorical components. Specifically, $z_{p,t}^k$ denotes the continuous observation of the k -th object’s pose, while $\mathcal{Z}_{s,t}$ and $\mathcal{Z}_{r,t}$ correspond to categorical observations associated with surfaces and rooms, respectively. The associated coverage terms, $v_{r,t}, v_{s,t} \in [0, 1]$, quantify the visibility of surface- and room-level observations. Further details are provided in a later section.

$P(x_{p,t}^k | x_{s,t}^k, x_{r,t}^k, z_{p,t}^k)$ denotes the probability that the k -th object is at a particular pose at time t , given that it is located on a specific surface, in a particular room, and the pose observation. $P(x_{s,t}^k | x_{r,t}^k, \mathcal{Z}_{s,t}, v_{s,t})$ represents the probability that the k -th object is located on a particular surface, given that it resides in a specific room and the corresponding surface observations. Lastly, $P(x_{r,t}^k | \mathcal{Z}_{r,t}, v_{r,t})$ represents the probability of the k -th object being in a particular room, given room observations. In this work, the first term is estimated using a particle filter, and the latter two are estimated using discrete Bayes filters.

Due to the similarity in their derivations, we present the full equations related to room-level belief updates $\text{bel}(x_{r,t}^k)$ in the main text, and omit the analogous equations for surface-level beliefs $\text{bel}(x_{s,t}^k)$.

1) *Estimating object’s semantic location:* For the k -th object, we apply hierarchical discrete Bayes filters to estimate its surface and room locations, respectively. Since we consider only stationary objects, typical motion models in recursive Bayesian filtering are unnecessary. For notational convenience, we use $\overline{\text{bel}}$ to denote the predicted belief in the following Bayes filtering equations, which is essentially the

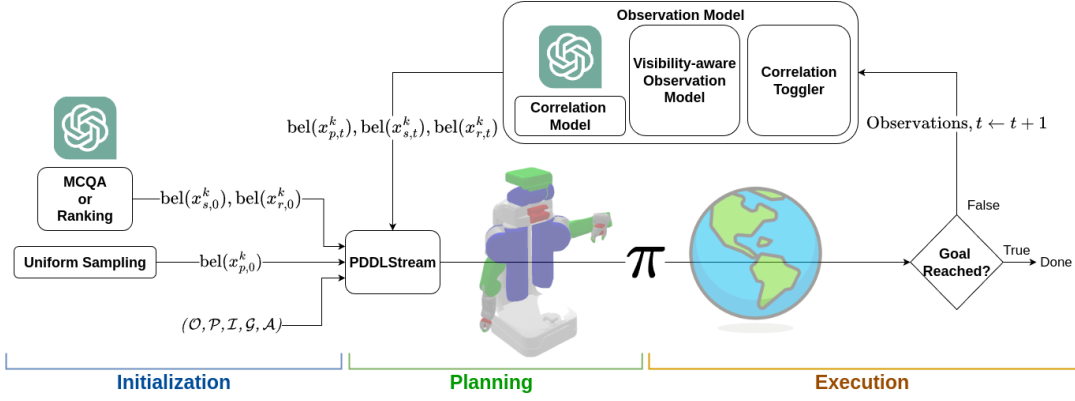


Fig. 1: The initial beliefs about the semantic locations of objects, $\text{bel}(x_{r,0}^k)$ and $\text{bel}(x_{s,0}^k)$, are derived from LLMs, while the initial beliefs about their poses, $\text{bel}(x_{p,0}^k)$, are uniformly distributed across all surfaces. The TAMP problem specification $(\mathcal{O}, \mathcal{P}, \mathcal{I}, \mathcal{G}, \mathcal{A})$, where the beliefs are incorporated into the cost of the observation action, is provided to a TAMP planner PDDLStream. The planner outputs a plan, which the robot executes. Upon executing the observation action, the co-location toggler determines whether to use the co-location model or not, based on the observed object. Then, the beliefs are updated using the proposed observation model, and the planner replans with the updated beliefs. Planning and execution are repeated until the goal state \mathcal{G} is reached. Further implementation details, derivations, additional experimental results, and link to the code are available at our project page: <https://coco-tamp.github.io>.

same as the prior belief at time $t - 1$. The recursive Bayes filter equation for room category estimation is as follows:

$$\text{bel}(x_{r,t}^k) = P(x_{r,t}^k | \mathcal{Z}_{r,t}, v_{r,t}), \quad (1)$$

$$= \eta P(\mathcal{Z}_{r,t} | x_{r,t}^k, v_{r,t}) \overline{\text{bel}}(x_{r,t}^k), \quad (2)$$

where η represents a normalization constant.

We denote $v_{r,t}, v_{s,t} \in [0, 1]$ as the visibility of a room and a surface, respectively, from the robot's field of view. These visibility values are computed using the particles sampled to estimate the posterior distribution $P(x_{p,t}^k | x_{s,t}^k, x_{r,t}^k, z_{p,t}^k)$. Let \mathcal{S} be the set of surfaces in a room. For each surface $s \in \mathcal{S}$, let $n_{s,0}$ be the number of particles initially drawn uniformly on s , and let $n_{s,t}^{(\text{seen})}$ be the number of particles detected on s at time t from robot's field of view when applying the detect action. Then, $v_{s,t}$ and $v_{r,t}$ are calculated as:

$$v_{s,t} = \frac{n_{s,t}^{(\text{seen})}}{n_{s,0}}, v_{r,t} = \frac{\sum_{s \in \mathcal{S}} n_{s,t}^{(\text{seen})}}{\sum_{s \in \mathcal{S}} n_{s,0}}. \quad (3)$$

The visibility is incorporated into our observation model to account for partial observability, where occlusions or limited field of view may prevent successful observations. Consequently, when reasoning about an object's semantic location (e.g., a specific surface within a room), a failed observation does not definitively indicate the object's absence.

Following COS-POMDP [13], the categorical observation set for a set of objects over semantic locations can be factored into a product of individual observations for each object using the conditional independence assumption [14], effectively separating the visibility-aware observation model for the object of interest during the process.

The conditional independence assumption suggests that the observation of an object z_t^i at the current time t is conditionally independent of all other objects' observation, given the object's state x_t^i . For example, given two objects

$i, j \in K$ where $i \neq j$,

$$P(z_t^i, z_t^j | x_t^j) = P(z_t^i | x_t^j) P(z_t^j | x_t^j). \quad (4)$$

Essentially, this assumption implies that the observation of object i does not provide additional information about the observation of object j once the state of object j is already known. Under this assumption, the observation model for updating $\text{bel}(x_{r,t}^k)$ can be written as:

$$P(\mathcal{Z}_{r,t} | x_{r,t}^k, v_{r,t}) = P(z_{r,t}^k | x_{r,t}^k, v_{r,t}) \times \prod_{j=1, j \neq k}^K P(z_{r,t}^j | x_{r,t}^k). \quad (5)$$

We refer to $P(z_{r,t}^k | x_{r,t}^k, v_{r,t})$ as *visibility-aware observation models*. The visibility-aware observation model is defined as follows:

- **Case 1:** Observation of the location where the object is placed (i.e., $z_{r,t}^k = x_{r,t}^k$).

$$P(z_{r,t}^k | x_{r,t}^k, v_{r,t}) = \begin{cases} (1 - p_{\text{fn}}) v_{r,t}, & \text{Det.}, \\ (1 - v_{r,t}) + v_{r,t} p_{\text{fn}}, & \text{No Det.} \end{cases} \quad (6)$$

- **Case 2:** Observation of the location where the object is not placed (i.e., $z_{r,t}^k \neq x_{r,t}^k$).

$$P(z_{r,t}^k | x_{r,t}^k, v_{r,t}) = \begin{cases} p_{\text{fp}} v_{r,t}, & \text{Det.}, \\ 1 - v_{r,t} p_{\text{fp}}, & \text{No Det.} \end{cases} \quad (7)$$

The false negative and false positive are denoted as p_{fn} and p_{fp} , respectively, which are set to 0.01 in our experiments.

2) *co-location model*: The visibility-aware observation model itself does not provide any information about the semantic location of another object. However, in real-world environments, similar objects are often stored together, while dissimilar objects are typically stored apart. To address this

gap, we propose a co-location model that leverages LLM embeddings.

Let \mathcal{R} be a set of rooms in an environment. Using the law of total probability, $P(z_{r,t}^j | x_{r,t}^k)$ can be rewritten as:

$$P(z_{r,t}^j | x_{r,t}^k) = \sum_{x_{r,t}^j \in \mathcal{R}} P(z_{r,t}^j | x_{r,t}^j) P(x_{r,t}^j | x_{r,t}^k). \quad (8)$$

We refer to $P(x_{r,t}^j | x_{r,t}^k)$ as *co-location models*, which incorporate the similarity between objects j and k .

We compute the similarity between objects by applying cosine similarity to their respective LLM embeddings. Specifically, for each object, we use an LLM to generate three sentences describing its common uses, concatenate these sentences into a single textual prompt, and then obtain an embedding of that prompt from an LLM embedding model. See Appendix VI-A for a prompting example. The cosine similarity between the resulting embeddings serves as the similarity score $\text{sim}(j, k) \in [-1, 1]$.

The co-location models exhibit the following properties:

- If objects j, k are identical ($\text{sim}(j, k) = 1$), the co-location models reduce to the Kronecker delta.
- If objects j, k have no co-location ($\text{sim}(j, k) = 0$), the co-location models become a uniform distribution.
- If objects j, k are very different ($\text{sim}(j, k) = -1$), the co-location models are the complement of the Kronecker delta.

If $\text{sim}(j, k) \geq 0$, the co-location models interpolate between the uniform distribution ($u(a, b)$) and the Kronecker delta ($\delta(a, b)$), depending on the similarity between the two objects. If $\text{sim}(j, k) \leq 0$, the co-location models interpolate between the normalized complement of the Kronecker delta and the uniform distribution, based on the similarity of the two objects. Formally, we define the co-location model $P(x_{r,t}^j | x_{r,t}^k)$ for R rooms in \mathcal{R} as shown below.

Let $\delta_{jk} = \delta(x_{r,t}^j, x_{r,t}^k)$, $\bar{\delta}_{jk} = \frac{1 - \delta(x_{r,t}^j, x_{r,t}^k)}{R-1}$, and $u_{jk} = u(x_{r,t}^j, x_{r,t}^k)$. Then

$$P(x_{r,t}^j | x_{r,t}^k) = \begin{cases} \text{sim}(j, k) \delta_{jk} \\ + (1 - \text{sim}(j, k)) u_{jk}, & \text{sim}(j, k) \geq 0, \\ \text{abs}(\text{sim}(j, k)) \bar{\delta}_{jk} \\ + (1 + \text{sim}(j, k)) u_{jk}, & \text{sim}(j, k) < 0. \end{cases} \quad (9)$$

As a result, observing an object increases the belief of similar objects sharing that semantic location, and vice versa for dissimilar objects.

Some objects may tend to be distributed across different areas of the environment; for instance, observing a light switch in a bedroom does not imply that all light switches are located in the bedroom. To accommodate such variability, CoCo-TAMP leverages the LLM to decide whether to enable the co-location model based on the semantics of the observed object. An example prompt used for this decision process is provided in Appendix VI-B.

3) *Estimating object pose*: The final step in object state estimation is to compute the belief about the object’s pose, $\text{bel}(x_{p,0}^k)$. To update $\text{bel}(x_{p,0}^k)$, two types of continuous observation models are introduced: one for when the object is visible, and the other for when the object is not visible. These observation models are applied within the particle filter algorithm.

$$\text{bel}(x_{p,t}^k) = P(x_{p,t}^k | x_{s,t}^k, x_{r,t}^k z_{p,t}^k), \quad (10)$$

$$= \eta P(z_{p,t}^k | x_{p,t}^k, x_{s,t}^k, x_{r,t}^k) \overline{\text{bel}}(x_p(t)), \text{ or} \quad (11)$$

$$= \eta P(z_{p,t}^k | x_{p,t}^k, x_{s,t}^k, x_{r,t}^k) \overline{\text{bel}}(x_p(t)). \quad (12)$$

$P(z_{p,t}^k | x_{p,t}^k, x_{s,t}^k, x_{r,t}^k)$ is used when the object is visible and increases the weights of particles representing poses close to the object’s observed pose. Specifically, we employ a Gaussian measurement model of the form:

$$P(z_{p,t}^k | x_{p,t}^{[i],k}, x_{s,t}^k, x_{r,t}^k) = \frac{1}{\eta} \exp\left(-\frac{d^2}{2\sigma^2}\right), \quad (13)$$

where d is the Euclidean distance between the particle’s pose and the observed pose. The parameters σ and η represent the standard deviation of the measurement noise and a normalization constant, respectively.

$P(z_{p,t}^k | x_{p,t}^k, x_{s,t}^k, x_{r,t}^k)$ is used when the object is not visible. This function retains the weights of particles that are occluded or outside the visible region, while decreasing the weights of particles that are visible.

$$P(\overline{z}_{p,t}^k | x_{p,t}^{[i],k}, x_{s,t}^k, x_{r,t}^k) = \begin{cases} 0, & \text{Visible,} \\ \frac{1}{\eta} w_{t-1}^{k,[i]}, & \text{Invisible,} \end{cases} \quad (14)$$

where η is a normalization constant. Similar to the co-location model, when object k is not visible but object j was observed, the weights of the particles are set using the semantic similarity between objects k and j as shown below.

$$w_t^k = \frac{1 + \text{sim}(j, k)}{2} e^{-d/\lambda} + \frac{1 - \text{sim}(j, k)}{2} (1 - e^{-d/\lambda}), \quad (15)$$

where d is the distance between the detected object and the particle, λ is a hyperparameter that determines how fast the weight should increase/decrease with respect to the distance.

IV. EXPERIMENTS

We evaluate CoCo-TAMP in household environments through both simulation and real-world experiments, measuring cumulative planning and execution time as well as the number of replanning iterations. In simulations, we conduct large-scale experiments leveraging a large household dataset (*i.e.*, the Housekeep dataset [15]) to place objects in a common-sense manner, validating the scalable efficiency of CoCo-TAMP compared to baselines. In real-world experiments, we demonstrate CoCo-TAMP’s practical effectiveness by deploying it on Toyota’s human support robot (HSR) [16]. In Figure 5 we compared initial belief generation across multiple LLM models. As GPT-4o consistently outperformed the alternatives within our evaluation scope, all subsequent experiments were conducted using GPT-4o via API calls.

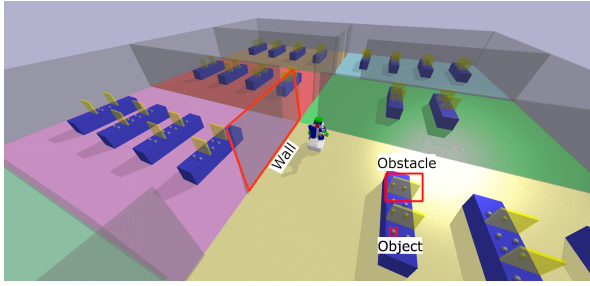


Fig. 2: Example of a simulated household environment.

A. Large-scale simulations

One challenge in generating simulated household environments is determining where to place objects, as our decisions may introduce bias in favor of our algorithm and fail to reflect common sense. To mitigate this issue, we sample household environments from the existing large-scale Housekeep dataset, which exhibits common-sense object placements.

An example of a simulated environment is shown in Fig. 2. Note that CoCo-TAMP does not use visual semantics and therefore a realistic simulation is unnecessary.

1) *Comparative analysis*: We evaluate the effectiveness of the proposed LLM-generated initial belief and co-location model by comparing six different variants. Our evaluation metrics include both the cumulative planning and execution time, as well as the number of replans necessary to solve the PO-TAMP problem with CoCo-TAMP. The six methods are: *Baseline*, *Co-Model*, *LLM generated belief update (LGBU)*, *MCQA*, *MCQA with Co-Model*, *MCQA with LGBU*.

These six methods differ in their method for generating the initial belief over the objects’ semantic locations (*i.e.*, $\text{bel}(x_{s,0}^k)$, $\text{bel}(x_{r,0}^k)$) as described in Section III-A, and the use of the co-location model as discussed in Section III-B.2. In particular, LGBU updates $\text{bel}(x_{s,t}^k)$ and $\text{bel}(x_{r,t}^k)$ directly using LLM-generated predictions based on the current observation using next-token prediction in a manner analogous to initial belief estimation (Section III-A). The specific prompts used for LGBU are detailed in Appendix VI-C. Table I summarizes these differences by illustrating both the initial belief generation method and the use of the co-location model or LLM-generated update for each method.

2) *Comparison results*: We design our experiments to test the following hypotheses:

- H1: CoCo-TAMP can generalize to diverse common household environments.
- H2: Common-sense knowledge about object locations, encoded in LLMs, can improve framework efficiency.
- H3: Semantic similarity between objects provides useful inductive bias that can improve framework efficiency.
- H4: Solely relying on LLMs for belief updates (LGBU) is insufficient for long-horizon planning and execution.

To test Hypothesis 1, we evaluate CoCo-TAMP in a large-scale simulated household environment across six layout configurations varying in the number of rooms and surfaces: 4 rooms and 8 surfaces, 4 rooms and 16 surfaces, 6 rooms and 12 surfaces, 6 rooms and 24 surfaces, 8 rooms and

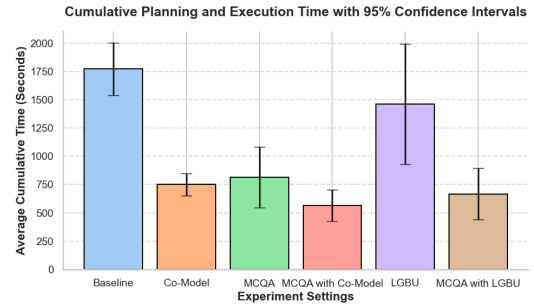


Fig. 3: Average cumulative planning and execution time over 50 environments in household layout with 6 rooms and 12 surfaces with 95% confidence interval.

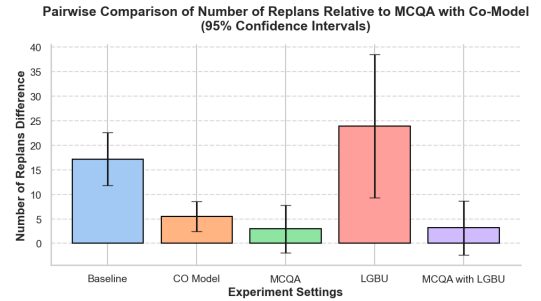


Fig. 4: Pairwise comparison of number of replans required to complete the task compared to MCQA with co-location Model with 95% confidence interval.

16 surfaces, and 8 rooms and 32 surfaces. For each layout, we sample 50 distinct environments by varying the labels of rooms, surfaces, and objects while keeping their quantities fixed. Occluding obstacles (*e.g.*, yellow boxes) are added to verify robustness under partial observability.

Given a task, such as detecting and moving an apple to the kitchen table, we record the cumulative planning and execution times and the number of replanning iterations across 50 environments for each layout. We also provide pairwise comparisons with respect to the best-performing method (*i.e.*, the one yielding the lowest cumulative planning and execution times and the fewest replanning iterations). These comparisons indicate, on average, the additional time and extra replanning iterations each method requires relative to the best-performing method. Notably, the best-performing method is deemed statistically superior when the mean value minus its confidence interval remains greater than zero. A comprehensive visualization of the results for layouts consisting of six rooms and twelve surfaces is presented in Figures 3 and 4, and the full results are provided in Table II and III.

To test Hypothesis 2, we compare the variant of CoCo-TAMP that uses only an LLM-generated initial belief (via MCQA) to a baseline without such semantic priors. The results show a notable reduction in both cumulative planning and execution time, as well as in the number of replanning iterations, indicating that the LLM prior effectively guides early decision-making.

To examine Hypothesis 3, a variant that uses only the

	Baseline	Co-Model	LGBU	MCQA	MCQA with Co-Model	MCQA with LGBU
$bel(x_{s,0}^k), bel(x_{r,0}^k)$	Uniform	Uniform	Uniform	MCQA	MCQA	MCQA
co-location model	No	Yes	No	No	Yes	No
LLM generated update	No	No	Yes	No	No	Yes

TABLE I: Comparison method details. The terms $bel(x_{s,0}^k)$ and $bel(x_{r,0}^k)$ denote the initial belief over surfaces and rooms, respectively, for object k .

	4 rooms 8 surfs	4 rooms 16 surfs	6 rooms 12 surfs	6 rooms 24 surfs	8 rooms 16 surfs	8 rooms 32 surfs
Baseline	632.5 ± 124.2	715.1 ± 240.5	1205.3 ± 272	1994 ± 602.7	1564.9 ± 362	3017.5 ± 1066.1
Co-Model	121.2 ± 58.4	70.2 ± 198.5	185.1 ± 145.6	355.3 ± 398.7	398.8 ± 272	0
MCQA	99.4 ± 95.3	169.9 ± 306.4	248.2 ± 262.2	1212.6 ± 735.1	134 ± 231	1925.3 ± 1320.8
MCQA with Co-Model	0	0	0	0	0	50.5 ± 834.8
LGBU	164.4 ± 92.9	605.8 ± 288.8	894.8 ± 551.9	509.2 ± 405	787.2 ± 321.3	1066.9 ± 1054.3
MCQA with LGBU	0.2 ± 48.8	339.2 ± 286.5	101.0 ± 210.7	351.1 ± 361.2	288.8 ± 308.3	408.5 ± 871.7

TABLE II: Pairwise comparison of cumulative planning and execution time (in seconds) relative to the best-performing method. Numbers represent the mean ± 95% confidence interval calculated over 50 problems, and 0 indicates the best-performing method.

	4 rooms 8 surfs	4 rooms 16 surfs	6 rooms 12 surfs	6 rooms 24 surfs	8 rooms 16 surfs	8 rooms 32 surfs
Baseline	9.5 ± 3.0	7.1 ± 6.8	17.2 ± 5.4	25.9 ± 8.8	19.8 ± 6.4	35.3 ± 13.7
Co-Model	3.0 ± 1.9	3.4 ± 5.8	5.5 ± 3.1	7.9 ± 5.5	10.6 ± 6.1	10.5 ± 7.4
MCQA	2.4 ± 2.7	0	2.9 ± 4.9	12.8 ± 8.7	0	13.9 ± 14.1
MCQA with Co-Model	0.6 ± 1.4	0.3 ± 6.6	0	0	1.4 ± 4.1	0
LGBU	4.1 ± 2.9	18.5 ± 10.7	23.9 ± 14.6	17.9 ± 10.1	21.3 ± 8.3	15.5 ± 8.8
MCQA with LGBU	11.2 ± 9.4	0	3.1 ± 5.5	10.2 ± 7.8	6.6 ± 7.0	6.1 ± 8.2

TABLE III: Pairwise comparison of number of replans relative to the best-performing method. Numbers represent the mean ± 95% confidence interval calculated over 50 problems, and 0 indicates the best-performing method.

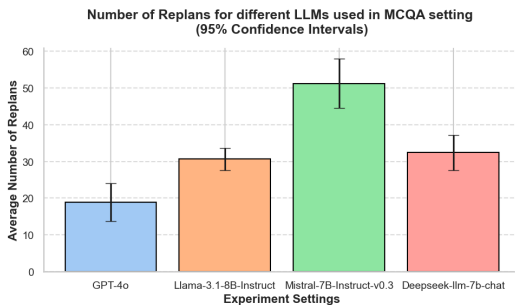


Fig. 5: Performance metrics over 50 environments in a household layout (6 rooms, 12 surfaces), comparing different LLMs for the MCQA setting. The plots show the number of replans with 95% confidence intervals. We see GPT-4o outperforms the other smaller models, and hence we have used GPT-4o in all other experiments in this paper.

co-location model during execution is compared against the baseline. This variant reduces the variability (confidence interval) as well as the average of cumulative time and number of replans. This supports that semantic co-locations between objects improve belief refinement throughout long-horizon planning and execution.

Finally, to investigate Hypothesis 4, we evaluate an approach referred to as LGBU, in which belief updates are performed solely using LLM-generated predictions. Experimental results indicate that LGBU results in more frequent replanning and higher cumulative planning and execution time compared to methods that utilize a principled Bayesian belief update, such as the Baseline and Co-Model variants. Figure 6 further reports results under adversarial configurations, where object placements were randomized to disrupt

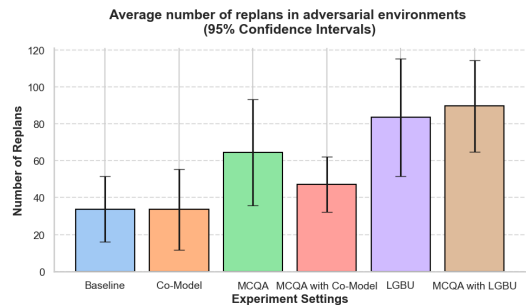


Fig. 6: Average number of replans over 5 environments in household layout with 6 rooms and 12 surfaces with 95% confidence interval. Note: The LGBU and MCQA with LGBU experiments had a total limit of 100 replans, which was exceeded 3 out of 5 times this experiment was run.

commonsense regularities. In this setting, Bayesian update methods consistently achieved task completion, whereas LGBU failed in 3 out of 5 trials. These findings indicate that while LLMs provide informative priors, ungrounded LLM-only updates lack the robustness required for reliable long-horizon execution.

Collectively, the results demonstrate that combining the LLM-generated initial belief with the co-location model yields the lowest average number of replanning iterations and cumulative planning and execution time, while also exhibiting the least variability across trials among all evaluated variants.

B. Real-world experiments

We demonstrate COCO-TAMP in a real-world setting using the HSR robot and a mock apartment environment

consisting of two rooms (living room, kitchen) and three surfaces (coffee table, bench, table).

In this setup, the coffee table and bench are located in the living room, while the table is in the kitchen. The task involves relocating an apple from the kitchen table to the coffee table. Additionally, we introduce a banana and a screwdriver to evaluate the effectiveness of the co-location model and use a cracker box and a cereal box to create occlusions. The real-world setup is shown in Figure 7.

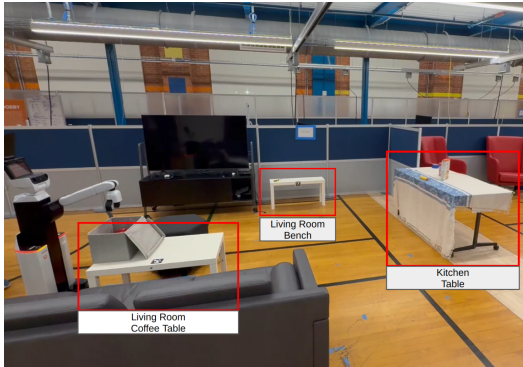


Fig. 7: Environment for real-world experiments.

The resulting cumulative planning and execution times are as follows: Baseline: 365 seconds, Co-Model: 112 seconds, and MCQA with Co-Model: 100 seconds.

V. RELATED WORK

Belief space planning aims to handle stochastic transitions and imperfect observations, requiring the solution of POMDP problems. Belief space TAMP planners [17], [18], [19], [20] typically focus on modeling uncertainty at both the task and motion levels while designing approximate solutions to mitigate the complexity of POMDPs. However, they inevitably require replanning to complete tasks within an integrated planning and execution framework. In contrast, CoCo-TAMP aims to enhance the efficiency of this belief space framework by shaping beliefs informed by LLMs.

As our work specifically addresses partial observability and even complete unobservability, the object search literature is also closely related [21], [13], [22], [23], [24]. While these methods tackle certain forms of POMDP problems, they do not fully address TAMP problems, as they do not incorporate task-level symbolic components. Therefore, the object search literature can be considered a subclass of CoCo-TAMP.

With the rise of LLM research, LLMs have been incorporated into TAMP problems in various ways, replacing components that are typically either engineered or learned from data—such as constraints [25], [26], [27], task-level plan computation [28], [29], [30], [31], [32], [33], or symbolic and geometric spatial relationships generator [34]—with their common-sense. Additionally, some approaches utilize LLMs to translate natural language descriptions of planning problems into PDDL representations [2], [35]. In contrast, CoCo-TAMP leverages LLMs to query a different aspect

of commonsense knowledge, specifically related to object placement.

VI. CONCLUSION

We present CoCo-TAMP, a belief-space planning and execution framework that leverages LLMs to integrate common-sense reasoning into PO-TAMP problems. The framework encodes two key types of common-sense knowledge: (1) certain objects are more likely to be found in specific locations, and (2) semantically similar objects tend to be co-located, while dissimilar ones are not. While LLMs likely possess transferable knowledge applicable to non-household domains (*e.g.*, factories or hospitals), CoCo-TAMP has not yet been evaluated in such settings. In some scenarios, full information about the environment layout may be unavailable, necessitating information gathering about rooms and surfaces. We leave these directions for future work, which we are currently pursuing.

APPENDIX

A. Co-location Model Prompt

Prompt: Explain the common use of a toaster in three sentences.

LLM Response (GPT4o): A toaster is commonly used to toast slices of bread, making them crispy and golden brown. It’s also useful for heating up and adding a slight crunch to bagels, English muffins, and some pastries. Additionally, many people use a toaster to quickly prepare breakfast items like Pop-Tarts or frozen waffles.

B. Co-Location Toggler Prompt

System Prompt: You will be given an object that commonly appears in a typical household environment. Using common sense, determine if the object tends to be distributed throughout a typical household, such as doorknobs and light switches.

User Prompt: The following is the object: [object]

LLM Response (GPT4o): True / False

C. LGBU Prompt

System Prompt: You will receive:

- Current belief about where an object might be located.
- `observation_loc`: the location that was just inspected.
- `visibility`: how much of the location was visible. 0 means not visible at all, 1 means fully visible.
- `result`: whether the object was found there.
- `co_detected`: other objects found at the same location.

Based on this information and common sense, predict where the object is most likely to be now. Choose the most likely location from the given options.

User Prompt:

- Current belief about [object_name]: [belief_text]
- `observation_loc`: [obs_loc]
- `visibility`: [visibility]
- `result`: [result]
- `co_detected`: [co_detected_text]

Given this information, where is object_name most likely to be? [options_text]

ACKNOWLEDGMENTS

This work has taken place in the Learning Agents Research Group (LARG) at UT Austin. LARG research is supported in part by NSF (FAIN-2019844, NRT-2125858), ONR (N00014-24-1-2550), ARO (W911NF-17-2-0181, W911NF-23-2-0004, W911NF-25-1-0065), DARPA(Cooperative Agreement HR00112520004 on Ad Hoc Teamwork) Lockheed Martin, and UT Austin’s Good Systems grand challenge. Peter Stone serves as the Chief Scientist of Sony AI and receives financial compensation for that role. The terms of this arrangement have been reviewed and approved by the University of Texas at Austin in accordance with its policy on objectivity in research.

REFERENCES

- [1] C. R. Garrett, R. Chitnis, R. Holladay, B. Kim, T. Silver, L. P. Kaelbling, and T. Lozano-Pérez, “Integrated task and motion planning,” *Annual review of control, robotics, and autonomous systems*, vol. 4, no. 1, pp. 265–293, 2021.
- [2] 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,” *arXiv preprint arXiv:2304.11477*, 2023.
- [3] K. Valmeekam, M. Marquez, S. Sreedharan, and S. Kambhampati, “On the planning abilities of large language models—a critical investigation,” *Advances in Neural Information Processing Systems*, vol. 36, pp. 75 993–76 005, 2023.
- [4] L. Guan, Y. Zhou, D. Liu, Y. Zha, H. B. Amor, and S. Kambhampati, “Task success is not enough: Investigating the use of video-language models as behavior critics for catching undesirable agent behaviors,” *arXiv preprint arXiv:2402.04210*, 2024.
- [5] S. Kambhampati, K. Valmeekam, L. Guan, M. Verma, K. Stechly, S. Bhambri, L. P. Saldyt, and A. B. Murthy, “Position: Llms can’t plan, but can help planning in llm-modulo frameworks,” in *Forty-first International Conference on Machine Learning*, 2024.
- [6] N. Hughes, Y. Chang, and L. Carlone, “Hydra: A real-time spatial perception system for 3d scene graph construction and optimization,” *arXiv preprint arXiv:2201.13360*, 2022.
- [7] C. R. Garrett, C. Paxton, T. Lozano-Pérez, L. P. Kaelbling, and D. Fox, “Online replanning in belief space for partially observable task and motion problems,” in *2020 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2020, pp. 5678–5684.
- [8] D. P. Bertsekas and J. N. Tsitsiklis, “An analysis of stochastic shortest path problems,” *Mathematics of Operations Research*, vol. 16, no. 3, pp. 580–595, 1991.
- [9] C. R. Garrett, T. Lozano-Pérez, and L. P. Kaelbling, “Pddlstream: Integrating symbolic planners and blackbox samplers via optimistic adaptive planning,” in *Proceedings of the international conference on automated planning and scheduling*, vol. 30, 2020, pp. 440–448.
- [10] B. Kumar, C. Lu, G. Gupta, A. Palepu, D. Bellamy, R. Raskar, and A. Beam, “Conformal prediction with large language models for multi-choice question answering,” *arXiv preprint arXiv:2305.18404*, 2023.
- [11] H. Vishwakarma, A. Mishler, T. Cook, N. Dalmaso, N. Raman, and S. Ganesh, “Improving decision-making in open-world agents with conformal prediction and monty hall,” in *NeurIPS 2024 Workshop on Open-World Agents*, 2024.
- [12] J. Su, J. Luo, H. Wang, and L. Cheng, “Api is enough: Conformal prediction for large language models without logit-access,” *arXiv preprint arXiv:2403.01216*, 2024.
- [13] K. Zheng, R. Chitnis, Y. Sung, G. Konidaris, and S. Tellex, “Towards optimal correlational object search,” in *2022 International Conference on Robotics and Automation (ICRA)*. IEEE, 2022, pp. 7313–7319.
- [14] A. P. Dawid, “Conditional independence in statistical theory,” *Journal of the Royal Statistical Society Series B: Statistical Methodology*, vol. 41, no. 1, pp. 1–15, 1979.
- [15] Y. Kant, A. Ramachandran, S. Yenamandra, I. Gilitschenski, D. Batra, A. Szot, and H. Agrawal, “Housekeep: Tidying virtual households using commonsense reasoning,” in *European Conference on Computer Vision*. Springer, 2022, pp. 355–373.
- [16] T. Yamamoto, K. Terada, A. Ochiai, F. Saito, Y. Asahara, and K. Murase, “Development of human support robot as the research platform of a domestic mobile manipulator,” *ROBOMECH journal*, vol. 6, no. 1, pp. 1–15, 2019.
- [17] L. P. Kaelbling and T. Lozano-Pérez, “Integrated task and motion planning in belief space,” *The International Journal of Robotics Research*, vol. 32, no. 9–10, pp. 1194–1227, 2013.
- [18] X. Fang, C. R. Garrett, C. Eppner, T. Lozano-Pérez, L. P. Kaelbling, and D. Fox, “Dimsam: Diffusion models as samplers for task and motion planning under partial observability,” in *2024 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2024, pp. 1412–1419.
- [19] A. Curtis, G. Matheos, N. Gothoskar, V. Mansinghka, J. Tenenbaum, T. Lozano-Pérez, and L. P. Kaelbling, “Partially observable task and motion planning with uncertainty and risk awareness,” *arXiv preprint arXiv:2403.10454*, 2024.
- [20] T. Pan, R. Shome, and L. E. Kavraki, “Task and motion planning for execution in the real,” *arXiv preprint arXiv:2406.03641*, 2024.
- [21] K. Zheng, Y. Sung, G. Konidaris, and S. Tellex, “Multi-resolution pomdp planning for multi-object search in 3d,” in *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2021, pp. 2022–2029.
- [22] J. Park, T. Yoon, J. Hong, Y. Yu, M. Pan, and S. Choi, “Zero-shot active visual search (zavis): Intelligent object search for robotic assistants,” in *2023 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2023, pp. 2004–2010.
- [23] S. Sharma, H. Huang, K. Shivakumar, L. Y. Chen, R. Hoque, B. Ichter, and K. Goldberg, “Semantic mechanical search with large vision and language models,” in *Conference on Robot Learning*. PMLR, 2023, pp. 971–1005.
- [24] K. Zheng, A. Paul, and S. Tellex, “Asystem for generalized 3d multi-object search,” in *2023 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2023, pp. 1638–1644.
- [25] A. Curtis, N. Kumar, J. Cao, T. Lozano-Pérez, and L. P. Kaelbling, “Trust the proc3s: Solving long-horizon robotics problems with llms and constraint satisfaction,” *arXiv preprint arXiv:2406.05572*, 2024.
- [26] W. Guo, Z. Kingston, and L. E. Kavraki, “Castl: Constraints as specifications through llm translation for long-horizon task and motion planning,” *arXiv preprint arXiv:2410.22225*, 2024.
- [27] N. Kumar, F. Ramos, D. Fox, and C. R. Garrett, “Open-world task and motion planning via vision-language model inferred constraints,” *arXiv preprint arXiv:2411.08253*, 2024.
- [28] S. Wang, M. Han, Z. Jiao, Z. Zhang, Y. N. Wu, S.-C. Zhu, and H. Liu, “Llm³: Large language model-based task and motion planning with motion failure reasoning,” *arXiv preprint arXiv:2403.11552*, 2024.
- [29] Z. Li, K. Yu, S. Cheng, and D. Xu, “League++: Empowering continual robot learning through guided skill acquisition with large language models,” in *ICLR 2024 Workshop on Large Language Model (LLM) Agents*, 2024.
- [30] A. Z. Ren, A. Dixit, A. Bodrova, S. Singh, S. Tu, N. Brown, P. Xu, L. Takayama, F. Xia, J. Varley *et al.*, “Robots that ask for help: Uncertainty alignment for large language model planners,” *arXiv preprint arXiv:2307.01928*, 2023.
- [31] A. Brohan, Y. Chebotar, C. Finn, K. Hausman, A. Herzog, D. Ho, J. Ibarz, A. Irpan, E. Jang, R. Julian *et al.*, “Do as i can, not as i say: Grounding language in robotic affordances,” in *Conference on robot learning*. PMLR, 2023, pp. 287–318.
- [32] W. Huang, P. Abbeel, D. Pathak, and I. Mordatch, “Language models as zero-shot planners: Extracting actionable knowledge for embodied agents,” in *International conference on machine learning*. PMLR, 2022, pp. 9118–9147.
- [33] M. Mohammadi, D. Honerkamp, M. Büchner, M. Cassinelli, T. Welschehold, F. Despinoy, I. Gilitschenski, and A. Valada, “More: Mobile manipulation rearrangement through grounded language reasoning,” *arXiv preprint arXiv:2505.03035*, 2025.
- [34] Y. Ding, X. Zhang, C. Paxton, and S. Zhang, “Task and motion planning with large language models for object rearrangement,” in *2023 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2023, pp. 2086–2092.
- [35] Y. Xie, C. Yu, T. Zhu, J. Bai, Z. Gong, and H. Soh, “Translating natural language to planning goals with large-language models,” *arXiv preprint arXiv:2302.05128*, 2023.