

Pragmatic Instruction Following and Goal Assistance via Cooperative Language-Guided Inverse Planning

Tan Zhi-Xuan*

Massachusetts Institute of Technology
xuan@mit.edu

Vikash Mansinghka

Massachusetts Institute of Technology
vkm@mit.edu

Lance Ying*

Harvard University
lanceying@seas.harvard.edu

Joshua B. Tenenbaum

Massachusetts Institute of Technology
jbt@mit.edu

ABSTRACT

People often give instructions whose meaning is ambiguous without further context, expecting that their actions or goals will disambiguate their intentions. How can we build assistive agents that follow such instructions in a flexible, context-sensitive manner? This paper introduces *cooperative language-guided inverse plan search* (CLIPS), a Bayesian agent architecture for pragmatic instruction following and goal assistance. Our agent assists a human by modeling them as a cooperative planner who communicates *joint plans* to the assistant, then performs *multimodal* Bayesian inference over the human’s goal from actions and language, using large language models (LLMs) to evaluate the likelihood of an instruction given a hypothesized plan. Given this posterior, our assistant acts to minimize expected goal achievement cost, enabling it to pragmatically follow ambiguous instructions and provide effective assistance even when uncertain about the goal. We evaluate these capabilities in two cooperative planning domains (Doors, Keys & Gems and VirtualHome), finding that CLIPS significantly outperforms GPT-4V, LLM-based literal instruction following and unimodal inverse planning in both accuracy and helpfulness, while closely matching the inferences and assistive judgments provided by human raters.

KEYWORDS

Inverse Planning; Bayesian Theory-of-Mind; Instruction Following; Human-Robot Cooperation; Value Alignment

ACM Reference Format:

Tan Zhi-Xuan*, Lance Ying*, Vikash Mansinghka, and Joshua B. Tenenbaum. 2024. Pragmatic Instruction Following and Goal Assistance via Cooperative Language-Guided Inverse Planning. In *Proc. of the 23rd International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2024)*, Auckland, New Zealand, May 6 – 10, 2024, IFAAMAS, 18 pages.

1 INTRODUCTION

Humans act upon the world through our words. We make requests, give instructions, and communicate information, so that we can better coordinate with each other [4]. In doing so, we are often

parsimonious, exploiting context to convey our intentions [24]. For example, if someone says “Can you hold that door?”, it is typically obvious which door they mean, even though nothing in the sentence distinguishes the door they just walked past from the one they are heading towards. We understand these requests because we interpret them *pragmatically*, in light of the goals and actions of others. How might we build assistive machines that do the same?

In this paper, we introduce *cooperative language-guided inverse plan search* (CLIPS), a Bayesian architecture for pragmatic instruction following and goal assistance (Figure 1). Building upon prior work in inverse planning [8, 57, 80], rational speech act theory [19, 23, 67], joint intentionality [66, 70, 78], assistance games [18, 25], and reward learning [34, 44, 56], CLIPS models humans as cooperative planners who communicate *joint plans* as instructions. Given this model, CLIPS performs *multimodal* goal inference from human actions and instructions, using a large language model (LLM) to score the likelihood of observed utterances [81], and computing a distribution over goals via sequential inverse planning [86]. This distribution then informs an assistive policy, which selects helpful actions under uncertainty about the human’s goal.

We evaluate CLIPS on a suite of multi-step goal assistance problems in a doors-and-keys gridworld [86] and the VirtualHome domain [54, 55]. In these problems, the assistant must infer the human’s goal from their actions and possibly ambiguous instructions, then decide how best to help. Even when leveraging LLMs, standard instruction following methods struggle with this setting because they disregard pragmatic context [2, 64, 75], while action-only goal inference [8, 55, 86] ignores linguistic information. Multimodal LLMs have access to all information, but they fail to ground it in a coherent theory-of-mind [12, 35]. In contrast, CLIPS is able to use observed actions and inferred goals to resolve *ambiguous language*, interpret *joint instructions*, and correct for *incomplete commands*, achieving much higher goal accuracy and cooperative efficiency than GPT-4V, LLM-based literal instruction following, and unimodal inverse planning, while correlating strongly with goal inference and assistance judgments provided by human raters.

2 COOPERATIVE LANGUAGE-GUIDED INVERSE PLAN SEARCH

We formulate the setting for CLIPS as a *language-augmented goal assistance game*, an extension of assistance games [18, 25] with linguistic utterances and uncertainty over goals (i.e. desired terminal states) [59]. We define this as a two-player Markov game between a human principal and an assistive agent, described by

*Equal contribution.



This work is licensed under a Creative Commons Attribution International 4.0 License.

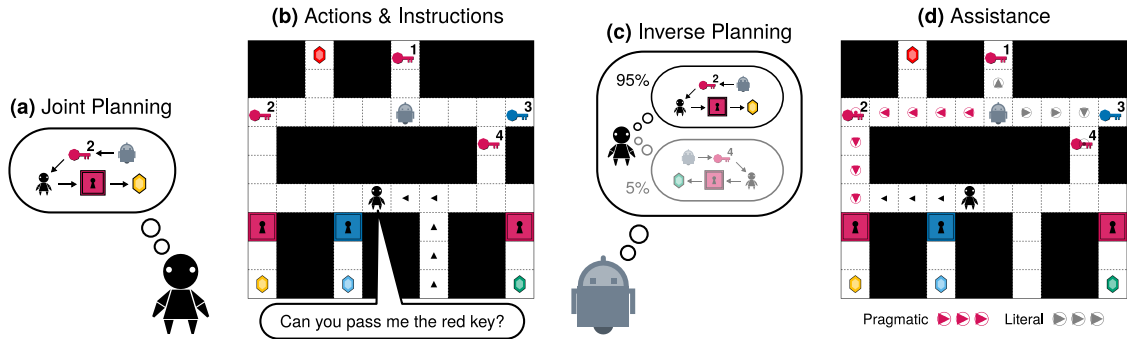


Figure 1: Overview of cooperative language-guided inverse plan search (CLIPS). We model a human principal as (a) cooperatively planning a joint policy for the human and the (robot) assistant. The human is (b) assumed to take actions from this joint policy while communicating planned actions as an instruction (“Can you pass me the red key?”). Observing this, (c) CLIPS infers the human’s goal and policy via Bayesian inverse planning. CLIPS then (d) acts by minimizing expected goal achievement cost, pragmatically interpreting the ambiguous instruction by picking up *Key 2*. In contrast, a literal instruction follower might pick up *Key 1* or *Key 4*, which are also red in color.

the tuple $(S, U, A^h, A^r, C, G, H, P_s, P_i)$, where S is the set of environment states, U a space of utterances u that the human may use to communicate at any step, A^h the set of human actions a^h , A^r the set of assistant actions a^r , C a set of cost functions $C : S \times U \times A^h \times A^r \rightarrow \mathbb{R}$ that map state-action transitions to real numbers, $G \subseteq \mathcal{P}(S)$ a set of possible goals g where each $g \subseteq S$ is a set of (terminal) states, H a horizon after which the game automatically terminates, $P_s(s'|s, a^r, a^h)$ the environment transition distribution, and $P_i(s_i, C, g)$ a distribution over initial states of the game. As in regular assistance games, the human knows (i.e. observes) the *true* cost function C and goal g sampled from P_i , but the assistant only observes the initial environment s_i . Thus, the assistant has to *infer* the true goal g .

Assistance games admit cooperative equilibria, but they are intractable to compute in general [18], and also assume more game-theoretic optimality from humans than may be warranted [49]. As such, our goal is to *approximately* solve the assistance game under reasonable modeling assumptions about how humans cooperate and communicate. By fixing a human model, the game becomes a partially observable Markov decision process (POMDP) from the assistant’s perspective [17], which can be solved through approximate methods [47, 52]. We describe this human model, then show how an assistant can perform Bayesian inference over such a model, using the acquired information to better assist humans.

2.1 Modeling cooperative action and communication

When humans cooperate, we direct our actions towards a shared goal while expecting that others will do the same. This capacity for *joint intentionality* [73] allows us to cooperate rapidly and flexibly while avoiding nested reasoning about each other’s minds [70, 78]. In CLIPS, we exploit this aspect of human cognition by modeling the human principal as a *cooperative planner*, who computes a *joint policy* π for both agents to achieve the goal g . The human follows π by taking an action a_t^h at each step t , and assumes the assistant will take an action a_t^r accordingly. These actions lead to a change in the state s_t . The human may also speak at any time t (with the

decision denoted by d_t), issuing a command c_t that summarizes their intended policy π . This command is then rendered in natural language as an utterance u_t . The overall generative process is shown in Figures 2a–b, and can be summarized as follows:

$$\text{Goal Prior:} \quad g \sim P(g) \quad (1)$$

$$\text{Joint Planning:} \quad \pi \sim P(\pi|g) \quad (2)$$

$$\text{Action Selection:} \quad a_t^h, a_t^r \sim P(a_t^h, a_t^r | s_t, \pi) \quad (3)$$

$$\text{Utterance Generation:} \quad u_t, c_t, d_t \sim P(u_t, c_t, d_t | s_t, \pi) \quad (4)$$

$$\text{State Transition:} \quad s_{t+1} \sim P(s_{t+1} | s_t, a_t^h, a_t^r) \quad (5)$$

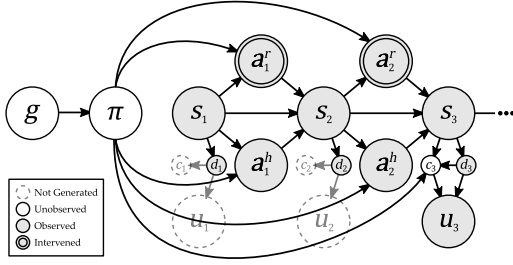
$P(g)$ is the assistant’s prior over the human’s goal $g \in G$, which we take to be uniform. We can also fold in uncertainty over the human’s cost function $C \in \mathcal{C}$ into this prior, treating g as a *specification* that includes both the goal condition and action costs.

To model joint planning, we assume that the human computes a Boltzmann policy π for goal g , which defines an approximately rational distribution over joint actions a_t^h, a_t^r :

$$\pi(a_t^h, a_t^r | s_t) = \frac{\exp(-\beta \hat{Q}_\pi(s_t, a_t^h, a_t^r))}{\sum_{\tilde{a}_t^h, \tilde{a}_t^r} \exp(-\beta \hat{Q}_\pi(s_t, \tilde{a}_t^h, \tilde{a}_t^r))} \quad (6)$$

Here, β is a rationality parameter controlling the optimality of the policy (higher β is more optimal), over which we may place a prior $P(\beta)$, and $\hat{Q}_\pi(s_t, a_t^h, a_t^r)$ is an estimate of the expected cost of reaching the goal g by taking actions a_t^h, a_t^r at state s_t . To estimate \hat{Q}_π efficiently, we extend prior work in online goal inference [86], using model-based planning to compute policies on-the-fly. In particular, we use real-time heuristic search (RTHS) as an anytime planner [9, 38, 40], which estimates Q -values in a neighborhood around the current state s_t via search (RTHS-POLICY-UPDATE in Fig. 2b), while using previously estimated Q -values to guide future searches. Unlike earlier methods for inverse planning, this avoids the intractability of estimating Q for all states and actions [56, 87].¹

¹As an additional simplification, we assume that the human and the assistant take turns while acting (i.e. a_t^r is a no-op when a_t^h is not, and vice versa). This reduces the branching factor while planning, but preserves the optimal solution [10].



(a) One graphical realization of our Bayesian model.

model CLIPS-MODEL(s_1, G, T)

```

 $g \sim \text{GOAL-PRIOR}(G)$ 
 $\pi \sim \text{POLICY-INIT}(g, s_1)$ 
for  $t \in [1, T]$  do
   $\pi \leftarrow \text{RTHS-POLICY-UPDATE}(\pi, g, s_t)$ 
   $a_t^h, a_t^r \sim \text{BOLTZMANN-DIST}(s_t, \pi)$ 
   $u_t, c_t, d_t \sim \text{UTTERANCE-MODEL}(s_t, \pi)$ 
   $s_{t+1} \sim \text{TRANSITION}(s_t, a_t^h, a_t^r)$ 
end for
end model

```

(b) CLIPS agent-environment model as a probabilistic program.

model UTTERANCE-MODEL(s_t, π)

parameters: $p_{\text{speak}}, H, K, \mathcal{E}$

$d_t \sim \text{BERNOULLI}(p_{\text{speak}})$

if $d_t = \text{TRUE}$ **then**

$a_{t:t+H} \leftarrow \text{ROLLOUT-POLICY}(s_t, \pi, H)$

$\alpha_{1:m} \leftarrow \text{EXTRACT-SALIENT-ACTIONS}(a_{t:t+H})$

$c_t \sim \text{RAND-SUBSET}(\alpha_{1:m}, K)$

$u_t \sim \text{LANGUAGE-MODEL}(c_t, \mathcal{E})$

else

$c_t, u_t \leftarrow \emptyset, ""$

end if

end model

(c) Utterance model $P(u_t, c_t, d_t | \pi)$ as a probabilistic program.

```

Command: (handover you me key1) where (iscolor key1 blue)
Utterance: Hand me the blue key.
Command: (unlock you key1 door1)
Utterance: Can you unlock the door for me?
Command: (pickup me key1) (unlock you key2 door1) where
(iscolor key1 blue) (iscolor door1 green)
Utterance: I'm getting the blue key, can you open the green door?
Command: (handover you me key1) (handover you me key2) where
(iscolor key1 green) (iscolor key2 red)
Utterance: Can you pass me the green and red keys?

```

(d) Paired examples \mathcal{E} of commands c_t and utterances u_t .

Figure 2: Model architecture. In CLIPS, we model the human as a cooperative planner who computes a joint policy π for a goal $g \in G$. The policy π dictates the human’s and assistant’s actions a_t^h, a_t^r at each state s_t , as well as the command c_t and utterance u_t that the human may decide d_t to communicate at step t . One realization of this process is depicted in (a), showing a case where an utterance u_3 is only made at $t = 3$. We implement this process as probabilistic program, shown in (b). Utterance generation is modeled by the subroutine in (c), which summarizes salient actions from policy π as a command c_t , then samples an utterance u_t using a (large) language model prompted with c_t and (d) a list of few-shot examples \mathcal{E} demonstrating how commands are translated into natural language.

Having computed the policy π , we model human communication as an approximately *pragmatic* process. Since the human wants to achieve their goal, they are likely to convey information about the policy they have in mind. However, long utterances are costly, and so the human is likely to convey only a salient summary of the policy, reducing the cost of communication. We capture these aspects of cooperative communication through a structured generation procedure that leverages the linguistic competence of LLMs:

- (1) At each step t , the human decides whether to communicate (d_t) with probability p_{speak} . If d_t is true, we rollout the policy π for H steps, producing a series of planned actions $a_{t:t+H}$.
- (2) The planned actions $a_{t:t+H}$ are filtered down to a domain-specific set of *salient* actions $\alpha_{1:m}$, and a random subset of up to K such actions are incorporated into a command c_t .
- (3) Finally, the command c_t is translated into a natural language utterance u_t using an LLM prompted with few-shot examples (Figure 2d). This defines a likelihood function $P(u_t | c_t)$.

This procedure is implemented as the probabilistic program shown in Figure 2c. While it does not capture all ways in which a pragmatic speaker might prefer some utterances over others, key features of pragmatic speech are modeled: Utterances are about plans, and hence practically useful, but are also restricted to small ($\leq K$) but salient subsets of those plans, ensuring they are informative without being costly. Furthermore, using an LLM as an utterance model $P(u_t | c_t)$ provides broad coverage, enabling CLIPS to handle much richer utterances than traditional pragmatic models [23].

2.2 Goal inference via inverse planning from actions and instructions

Using this probabilistic model, our assistant can infer a posterior distribution over the human’s goals g and policies π given a sequence of observed actions $a_{1:T}^h, a_{1:T}^r$, states $s_{1:T}$, utterance decisions $d_{1:T}$, and utterances $u_{1:T}$ (where u_t is the empty string whenever $d_t = \text{FALSE}$). Some care must be taken, however. Since the assistant is also taking actions $a_{1:T}^r$ alongside the human principal, it is tempting to infer the human’s goal by conditioning on *both* agents’ actions:

$$\begin{aligned}
 & P(g, \pi | s_{1:T}, d_{1:T}, u_{1:T}, a_{1:T}^h, a_{1:T}^r) \\
 & \propto P(g, \pi, s_{1:T}, d_{1:T}, u_{1:T}, a_{1:T}^h, a_{1:T}^r) \\
 & = P(g, \pi) \prod_{t=1}^T P(s_t | s_{t-1}, a_{t-1}^{hr}) P(u_t, d_t | s_t, \pi) \underbrace{P(a_t^h, a_t^r | s_t, \pi)}
 \end{aligned} \tag{7}$$

For an *external* observer, this is the appropriate distribution to compute [63, 81]. However, prior work on cooperative inference also computes Equation 7 when observers are *internal* to the environment [70, 78], which can lead to pathologies. For example, the assistant might condition on the fact that it has not moved so far ($a_{1:T}^r$ are no-ops), leading it to infer that the human is pursuing a goal g where the assistant’s help is unnecessary. To address this, it is crucial to recognize that the assistant is *intervening* upon the environment through its actions (Figure 2a, double-lined nodes).²

²Alternatively, it is enough to note the assistant’s actions are selected according to *different policy* π' than the inferred joint policy π computed by the human. Thus it is *safer* for the assistant not to update its beliefs as if its actions come from π .

As such, we should detach the evidential connection between the assistant’s actions $a_{1:T}^r$ and any causal ancestors, which we denote using Pearl’s **do**-operator [53]:

$$\begin{aligned} & P(g, \pi | s_{1:T}, d_{1:T}, u_{1:T}, a_{1:T}^h, \mathbf{do}(a_{1:T}^r)) \\ & \propto P(g, \pi, s_{1:T}, d_{1:T}, u_{1:T}, a_{1:T}^h, \mathbf{do}(a_{1:T}^r)) \quad (8) \\ & = P(g, \pi) \prod_{t=1}^T P(s_t | s_{t-1}, a_{t-1}^{hr}) P(u_t, d_t | s_t, \pi) \underline{P}(a_t^h | s_t, \pi) \end{aligned}$$

The difference between Equations 7 and 8 lies in the final underlined term: Whereas an external observer reweights its beliefs by incorporating the likelihood of both agents’ actions $P(a_t^h, a_t^r | s_t, \pi)$, our assistive agent should only incorporate the likelihood of the human’s action $P(a_t^h | s_t, \pi)$.

We compute the distribution in Equation 8 sequentially, as shown in Algorithm 1. An initial belief b_0 is returned by BELIEF-INIT, which generates a set of N weighted samples by either sampling or enumerating over the goal prior $P(g)$ and policy distribution $P(\pi|g)$ (accounting e.g. for uncertainty over the rationality parameter β). Then at each step t , the previous belief b_{t-1} is updated by adjusting the weight w^i associated with each goal g^i and policy π^i sample. To do this, BELIEF-UPDATE first refines the Q -value estimates that define the policy π^i by running more iterations of RTHS (RTHS-POLICY-UPDATE in L9). The weights are then updated to reflect the new state s_t and decision to speak d_t (L10-11). If d_t is true (i.e. u_t is observed), we also marginalize over all commands c_t that could have generated u_t from π , and update w^i with the resulting mixture likelihood (L13-14). Finally, we condition on the human’s action a_t^h , and return the resulting sample collection (L16-18).

If we enumerate over all possible goals and policy configurations in BELIEF-INIT, this is an *exact* Bayesian inference algorithm. In environments where there are too many hypotheses for this to be efficient, Algorithm 1 can readily be extended to a sequential Monte Carlo algorithm [16, 42] in the style of sequential inverse plan search [86]. For our experiments, however, enumeration is feasible, and we opt for this approach to avoid variance from sampling.

2.3 Pragmatic instruction following as goal assistance under uncertainty

With Algorithm 1, our assistive agent is able to infer a distribution over the human’s goal g and policy π at each step t . How should it use this information to act, especially if the goal g remains uncertain? In keeping with our Bayesian approach, we minimize *expected* cost, acting to help the human achieve their likely goal as quickly as possible [61]. Since each inferred policy π^i corresponds to a \hat{Q}_{π^i} -value function, the assistant can select actions by minimizing the expected \hat{Q}_{π^i} value given its uncertainty over π^i :

$$a_t^{r*} = \arg \min_{a_t^r} \mathbb{E}_{\pi^i} [\hat{Q}_{\pi^i}(s_t, a_t^h, a_t^r)] \quad (9)$$

We implement this assistance policy with Algorithm 2, which effectively solves the Q_{MDP} approximation of the assistive POMDP [27, 45, 47]. At each step t , we assume the assistant has already updated its belief b_{t-1} from the previous step $t-1$, and also observes the human’s action ³. We then iterate over all samples (w^i, g^i, π^i) , update the policies π^i if necessary, and perform a weighted sum of

³This assumption is natural in our turn-based setting. When actions are simultaneous the assistant can minimize expectation over both policies π^i and human actions a_t^h .

Algorithm 1 CLIPS belief initialization and update

```

1: procedure BELIEF-INIT( $N$ )
2:    $w^i \leftarrow 0$ ;  $g^i \sim P(g)$ ;  $\pi^i \sim P(\pi|g^i)$  for  $i \in [1, N]$ 
3:    $b_0 \leftarrow \{(w^i, g^i, \pi^i)\}_{i=1}^N$ ; return  $b_0$ 
4: end procedure
5:
6: procedure BELIEF-UPDATE( $b_{t-1}, s_{t-1:t}, d_t, u_t, a_{t-1:t}^h, a_{t-1}^r$ )
7:    $\{(w^i, g^i, \pi^i)\}_{i=1}^N \leftarrow b_{t-1}$ 
8:   for  $i \in [1, N]$  do
9:      $\pi^i \leftarrow \text{RTHS-POLICY-UPDATE}(\pi^i, g^i, s_t)$ 
10:     $w^i \leftarrow w^i \cdot P(s_t | s_{t-1}, a_{t-1}^h, a_{t-1}^r)$ 
11:     $w^i \leftarrow w^i \cdot P(d_t | s_t, \pi)$ 
12:    if  $d_t = \text{TRUE}$  then
13:       $P(u_t | s_t, \pi) \leftarrow \sum_{c_t} P(u_t | c_t) P(c_t | s_t, \pi)$ 
14:       $w^i \leftarrow w^i \cdot P(u_t | s_t, \pi)$ 
15:    end if
16:     $w^i \leftarrow w^i \cdot P(a_t^h | s_t, \pi)$ 
17:  end for
18:   $w^i \leftarrow w^i / \sum_{j=1}^N w_j$ 
19:   $b_t \leftarrow \{(w^i, g^i, \pi^i)\}_{i=1}^N$ ; return  $b_t$ 
20: end procedure

```

Algorithm 2 CLIPS Q_{MDP} assistance policy

```

1: procedure ASSISTANCE-POLICY( $b_{t-1}, s_t, a_t^h$ )
2:   for  $a_t^r \in \mathcal{A}^r(s_t)$  do
3:      $\hat{Q}_{\text{assist}}(s_t, a_t^h, a_t^r) \leftarrow 0$ 
4:     for  $(w^i, g^i, \pi^i) \in b_{t-1}$  do
5:        $\pi^i \leftarrow \text{RTHS-POLICY-UPDATE}(\pi^i, g^i, s_t)$ 
6:        $\hat{Q}_{\text{assist}}(s_t, a_t^h, a_t^r) \leftarrow w^i \cdot \hat{Q}_{\pi^i}(s_t, a_t^h, a_t^r)$ 
7:     end for
8:   end for
9:    $a_t^{r*} \leftarrow \arg \min_{a_t^r} \hat{Q}_{\text{assist}}(s_t, a_t^h, a_t^r)$ 
10:  return  $a_t^{r*}$ 
11: end procedure

```

\hat{Q}_{π^i} -values for all assistant actions a_t^r according to their inferred weights w^i . This produces a vector of *assistive* Q -values, which we minimize over to compute the best assistive action a_t^{r*} .

3 EXPERIMENTS

We evaluate our method against a variety of baselines in two domains: a cooperative gridworld puzzle called **multi-agent Doors, Keys & Gems** (mDKG) originally developed in [81, 86], and **VirtualHome** (VH), a virtual household simulator [54, 55]. In both domains, a human principal and a (robot) assistant cooperate on multi-step tasks to accomplish the human’s goals. In mDKG, the goals are four colored gems which are often secured behind doors. Doors can be unlocked by keys of the same color, and the assistant can help by collecting keys or unlocking doors (Figure 1). In VirtualHome, the goals are 6–12 household tasks, including setting up tables and preparing ingredients (Figure 3). Both domains admit encodings in the Planning Domain Definition Language (PDDL)



Figure 3: Example goal assistance problem in VirtualHome, where the principal and assistant collaborate to set the dinner table. The principal places three plates on the table, then says “Could you get the forks and knives?”. A pragmatic assistant has to infer the number of forks and knives from context (in this case, three each).

[48] with some extensions [84], which we use as the environment representation for our planning and inverse planning algorithms.

To systematically test the ability of assistive agents to pragmatically follow a human’s instructions and assist with their goals, we developed a dataset of *goal assistance problems* for each domain (30 problems in mDKG, 25 problems in VH). In each problem, the goal is initially unknown to the assistant. While taking a sequence of actions $a_{1:T}^h$, the human communicates one or more instructions $u_{1:T}$ to the assistant. The assistant then has to make an inference about the human’s goal g , and chooses actions to best assist them. We designed our problems to cover a range of scenarios inspired by human communication we observed in exploratory studies, varying the information that can be deciphered from the human’s actions or instructions (see Figure 4 and Supplementary Information). In mDKG, we restricted the human to only one instruction u_T , since more information would render goal inference trivial. In VH, however, the goal space was larger, allowing us to construct 10 problems with multiple utterances such as “Can you get the plates?” followed by “Bring the bowls too!”

3.1 Model configuration

In our experiments, we configured CLIPS to compute Q -values using the real-time adaptive A^* variant of RTHS [38], and a Gamma prior over the rationality parameter β . For the utterance model, we used either OpenAI’s 6.7B curie model [11] or the davinci-002 as our LLMs, due to the more diffuse probabilities they provided as base models. For goal assistance, we ran Algorithm 2 in offline mode, fixing the inferred posterior after observing $a_{1:T}^h$ and $u_{1:T}$, then selecting actions by minimizing expected Q -values with respect to the fixed posterior. This was evaluated against a simulated human agent which followed a joint policy to the true goal, unless it was apparent that the assistant was not doing the same. While Algorithm 2 can also be run in online mode, we fixed the posterior to better match the information that our human raters were provided.

3.2 Baselines

To evaluate the benefit of pragmatic instruction following over either a literal interpretation of instructions, or goal inference from only one modality, we included the following baselines:

3.2.1 Unimodal Inverse Planning. We implemented action-only and language-only inverse planning (IP) baselines as ablations of CLIPS, where the action-only baseline is similar to sequential inverse plan search [86] for inference and Watch-and-Help [55] for assistance. In these baselines, the assistive agent uses the same belief update shown in Algorithm 1, except that it conditions on only actions $a_{1:t}^h$ or only the utterance u_t . In cases with no actions, CLIPS and the language-only baseline are equivalent.

3.2.2 LLM-Based Literal Listener. The literal listener baseline interprets the instruction u_t in state s_t without further information about the human’s actions or goals. To implement this baseline, we adapt the utterance model in Figure 2c to sample from the space of all assistant-directed commands c_t that are possible in state s_t . In particular, we enumerate over all assistant actions \mathcal{A}^r , and select a subset of up to K salient actions to form a command c_t . This defines a distribution over commands $P(c_t|s_t)$ which *does not* depend the principal’s policy π or goal g . Given the command c_t , we use an LLM (text-davinci-002 for mDKG, davinci-002 for VH) as an utterance likelihood $P(u_t|c_t)$. Assuming a uniform prior over commands, we can perform enumerative inference to compute the posterior $P(c_t|u_t, s_t)$ given an instruction u_t in state s_t .

Given a command c_t , the assistant still needs to generate an assistive plan. We did this in two ways: The *naive* literal assistant interprets a command c_t by randomly selecting one concrete grounding (e.g. “Can you pick up the red key?” could mean picking up the red key closest to the human or some other key), then planning to achieve that concrete goal. In contrast, the *efficient* literal assistant tries to directly satisfy the command c_t in the most efficient way⁴ (e.g. “Could you unlock the blue door?” is satisfied by unlocking the blue door closest to the assistant). For both variants, the assistant can either satisfy the most likely command c_t^* , or sample a command from the full distribution $P(c_t|u_t, s_t)$. We report results for the latter after averaging over 10 samples, using systematic sampling to reduce variance.

3.2.3 Multimodal LLM (GPT-4V). For the mDKG domain, we used GPT-4 with Vision (GPT-4V) [1] as a purely neural baseline, prompting it with the same set of rules and instructions that we provided to our human raters, along with the final frame of each animated visual that we showed to humans (similar to Figure 1(b)). This allowed us to probe the degree to which multimodal LLMs are capable of intuitive pragmatic reasoning given spatially-grounded actions and verbal instructions [12]. Due to cost and rate limits, we report zero-shot performance with a temperature of zero.

More details about our models and datasets can be found in the Supplementary Information (<https://osf.io/v8ru7/>). Source code is available at <https://github.com/probcomp/CLIPS.jl>.

3.3 Human judgments

As an additional standard for comparison, we conducted a study with 100 human participants from the US through Prolific (mean age = 39.8, 59 men, 38 women, 2 non-binary), presenting each of them with 15 goal assistance problems from the mDKG domain. In each problem, participants saw the actions taken by the human

⁴Note that the naive/efficient distinction was not meaningful in the VirtualHome problems, where utterance ambiguity was not due to grounding ambiguity.

Table 1: Performance of CLIPS vs. baseline methods, measured in terms accuracy (posterior probability of true goal, precision and recall for assistance options), helpfulness (plan length and human cost relative to CLIPS), and human similarity (correlation of goal inferences and assistance options with mean human ratings) Metrics are averaged across the dataset per domain, with standard errors reported in brackets.

Method	$P(g_{\text{true}})$	Accuracy		Helpfulness		Human Similarity	
		Assist. Prec.	Assist. Rec.	Rel. Plan Length	Rel. Human Cost	Goal Cor.	Assist. Cor.
<i>Doors, Keys & Gems</i>							
Humans	0.67 (0.04)	0.83 (0.02)	0.85 (0.01)	–	–	–	–
CLIPS (Ours)	0.74 (0.05)	0.97 (0.03)	0.97 (0.03)	1.00 (0.00)	1.00 (0.00)	0.93 (0.01)	0.96 (0.01)
Lang. Only IP	0.55 (0.05)	0.90 (0.06)	0.83 (0.06)	1.26 (0.10)	1.18 (0.07)	0.74 (0.01)	0.83 (0.01)
Action Only IP	0.31 (0.04)	0.43 (0.09)	0.40 (0.09)	1.68 (0.15)	1.46 (0.08)	0.15 (0.01)	0.45 (0.01)
Literal Efficient	–	0.65 (0.08)	0.54 (0.07)	1.55 (0.11)	1.58 (0.10)	–	0.47 (0.01)
Literal Naive	–	0.58 (0.04)	0.47 (0.02)	1.54 (0.07)	1.52 (0.05)	–	0.54 (0.01)
GPT-4V	0.29 (0.08)	0.39 (0.08)	0.37 (0.08)	–	–	0.10 (0.01)	0.11 (0.01)
<i>VirtualHome</i>							
CLIPS (Ours)	0.63 (0.04)	0.87 (0.04)	1.00 (0.00)	1.00 (0.00)	1.00 (0.00)	–	–
Lang. Only IP	0.37 (0.05)	0.59 (0.05)	0.96 (0.03)	1.33 (0.07)	1.35 (0.08)	–	–
Action Only IP	0.25 (0.03)	0.61 (0.07)	0.84 (0.07)	1.30 (0.07)	1.33 (0.07)	–	–
Literal Listener	–	0.64 (0.06)	0.70 (0.07)	1.48 (0.08)	1.54 (0.09)	–	–

principal and observed the instruction at the end. Participants were then asked to select the set of gems that they thought were likely to be principal’s goal, and then indicate how the robot agent should best assist the principal. This gave 50 goal and assistance ratings for each problem (93% power for Cohen’s $d=0.5$ at $\alpha=0.05$).

For the assistance question, we presented each participant with a set of *assistance options*, corresponding to either picking up keys or unlocking doors. This allowed us to query participants’ understanding of how to provide assistance without having them solve the entire problem. To compare these human-provided assistance options with CLIPS and our baselines, we extracted the corresponding actions from the assistive plans generated by each method, and (where applicable) estimated the marginal probability of a particular assistance option occurring in the assistive plan via sampling.

3.4 Performance metrics

We evaluated CLIPS and the baseline methods in terms of their goal and assistance accuracy, the helpfulness and efficiency of the generated plans, and their similarity to human judgments. To evaluate accuracy, we calculated the probability $P(g_{\text{true}})$ assigned to the true goal for the non-literal methods, as well as precision and recall for selecting the optimal assistance options that were consistent with the instruction. To evaluate helpfulness and efficiency, we computed the length of the generated plan and the total action cost incurred by the human principal, relativized to CLIPS. Finally, to evaluate human similarity in mDKG, we calculated the correlation between each method’s outputs with the average ratings provided by our participants. We calculated Pearson’s r as well its variability from 1000 bootstrapped samples of the human dataset.

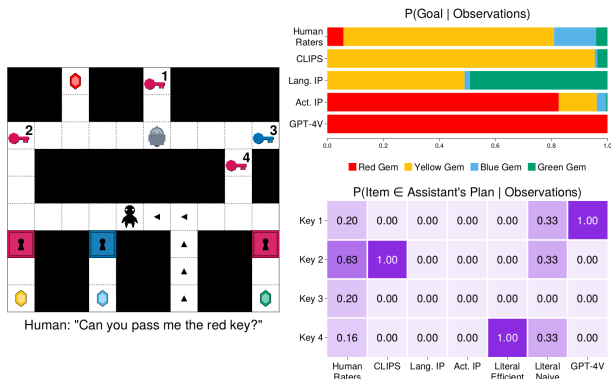
3.5 Results

The results of our experiments are presented in Table 1. In terms of accuracy, CLIPS assigned higher probability to the true goal than the baselines (note that 100% is not possible since goals cannot always be distinguished, as in Fig. 4d), while achieving close to perfect precision and recall in selecting the optimal assistance options.

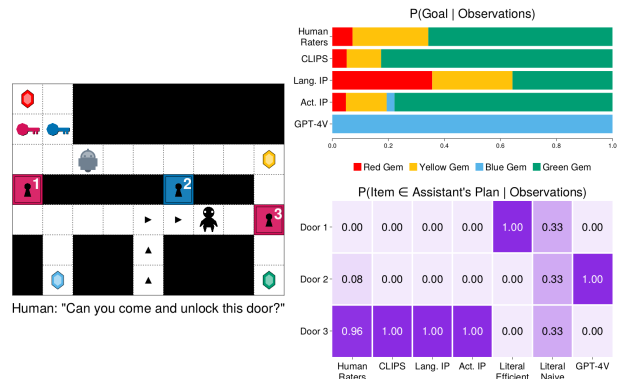
Indeed, CLIPS outperformed the average human in this regard. In contrast, unimodal inverse planning was much less accurate at inferring the true goal, and all baselines selected the correct assistance options at significantly lower rates. This affected the efficiency of assistance: CLIPS was 1.26 to 1.68 times faster in achieving the human’s goal than the baselines, and produced plans that were 1.18 to 1.58 times less costly for the human.

For human similarity in the mDKG domain, we found that both the goal inferences and the assistance options produced by CLIPS correlated highly with human ratings, achieving Pearson’s r of 0.93 (95% CI: 0.91–0.94) and 0.96 (95% CI: 0.95–0.96) respectively. In contrast, unimodal inverse planning produced goal inferences that were highly dissimilar from humans. Correlation with human-selected assistance options was also poor. These results demonstrate the importance of accounting for pragmatic context. Out of all methods, GPT-4V performed the worst, possibly because of the spatial reasoning and multi-step planning required for coherent goal inference in mDKG. Among participants themselves, we found that the median Pearson’s r of each participant’s ratings with the mean rating was 0.85 (IQR: 0.71–0.94) for goal inferences and 0.87 (IQR: 0.65–0.95) for assistance options, indicating that our dataset of human judgments is a reliable measure of average human performance.

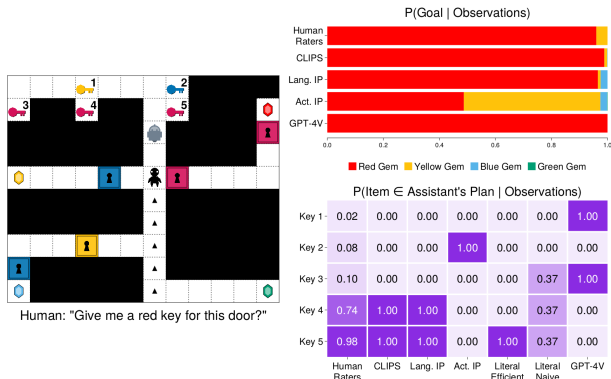
In Figure 4, we compare the results of CLIPS against human judgements and baselines on six illustrative goal assistance problems, providing qualitative observations in the captions. The results show how CLIPS closely mirrors human judgements for both goal inference and assistance, resolving ambiguity in predicates (Fig. 4a) and indexicals (Fig. 4b) while successfully completing partial instructions (Fig. 4c) and interpreting joint instructions (Fig. 4e). In comparison, the unimodal and literal baselines make less confident inferences or fail to assist appropriately, while GPT-4V provides incoherent answers. We also find that both humans and CLIPS are able to assist appropriately when there is significant goal uncertainty (Fig. 4d), and even when the optimal assistive plans for each goal make diverging recommendations (Fig. 4f), illustrating the importance of uncertainty-aware assistance.



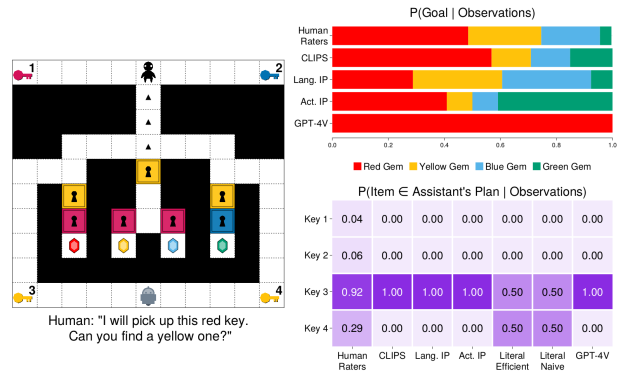
(a) **Ambiguous Predicates.** The human player asks for a red key when there are three available. CLIPS resolves this ambiguity, inferring that the human wants *Key 2* to reach the yellow gem.



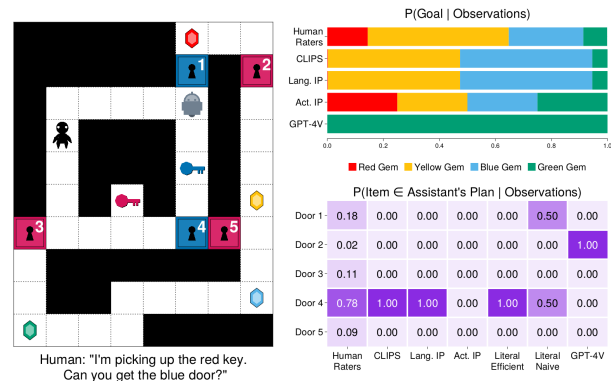
(b) **Ambiguous Indexicals.** The human player refers to a door using the indexical "this". CLIPS interprets this pragmatically, inferring that *Door 3* is intended even though *Door 2* is equally close to the human.



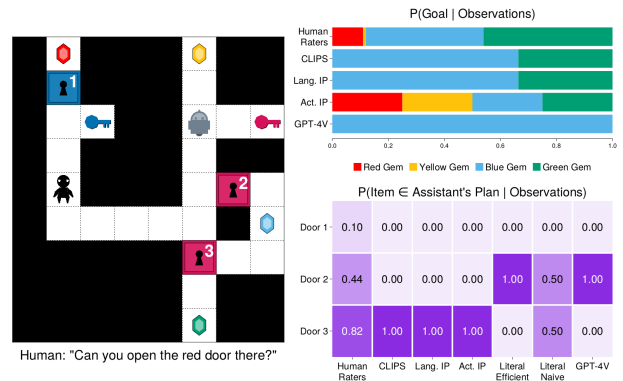
(c) **Partial Instructions.** The human player asks for a single red key, but requires two to reach the red gem. CLIPS infers this and picks up both keys (*Key 4* and *Key 5*), unlike the literal listener baselines.



(d) **Uncertain Goals.** The human player asks for a yellow key while going to pick up a red key. CLIPS pragmatically infers their goal to be one of the three gems on the left, and hence picks up the left key (*Key 3*).



(e) **Joint Instructions.** The human player asks for a blue door to be unlocked while indicating that they will pick up a red key. The CLIPS assistant infers that it should unlock *Door 4* so that the human can unlock *Door 5*.



(f) **Safe Assistance.** The human player asks the robot to unlock a red door when there are two such doors. The CLIPS assistant reasons that unlocking *Door 3* is safer as it leads to *both* blue and green gems, unlike *Door 2*.

Figure 4: Goal assistance problems in Doors, Keys & Gems. Each sub-figure contains a visual (*left*), instruction (*bottom left*), goal posteriors produced by each method (*top right*), and the probability of a key or door appearing in the assistive plans generated by each method (*bottom right*). Our pragmatic goal assistance method, CLIPS, best matches the goal inferences and assistance options produced by human raters (averaged across raters). In contrast, language and action-only inverse planning (Lang. IP & Act. IP) have higher goal uncertainty, the literal baselines fail to resolve instruction ambiguity, and GPT-4V often produces incoherent responses.

4 RELATED WORK

Theory of Mind as Bayesian Inverse Planning. Bayesian Theory of Mind (BToM) posits that humans reason about the actions and mental states of others through (approximately) Bayesian inference [6, 7]. In particular, Bayesian inverse planning can model how humans infer the goals of others by assuming that agents act rationally to achieve their goals [8, 21, 33, 63]. To efficiently solve these inference problems, prior work uses sequential Monte Carlo methods [16, 42] to perform online goal inference over model-based planning agents [3, 86]. CLIPS builds upon this paradigm, extending recent work by Ying et al. [81] on multi-agent inverse planning.

Rational Speech Acts. Rational Speech Act (RSA) theory frames communication as a recursive process, where speakers make pragmatic utterances that optimize relevance to listeners, and listeners interpret speaker utterances in light of the pragmatic goals that underlie those utterances [19, 23, 67]. CLIPS instantiates the RSA framework in the context of cooperative planning.

Multimodal Goal Inference and Reward Learning. Goal inference can be framed as online inverse reinforcement learning (IRL), where the aim is to learn a reward or cost function explaining the agent’s behavior in a single episode [32]. CLIPS can thus be viewed as a language-informed IRL method, though IRL is primarily applied offline given a dataset of expert demonstrations [20, 74, 77]. Particularly relevant is *reward-rational implicit choice* [34], a framework for multimodal Bayesian reward learning from heterogeneous human feedback. CLIPS can be seen as an application of this framework for specific modalities, but extended to the cooperative setting.

Decentralized Cooperation and Joint Intentionality. When multiple agents cooperate on a task, they usually need to track each other’s mental states to make decisions. Common approaches to this involve recursive reasoning over each other’s goals and plans, which can quickly grow intractable [13, 17]. To address this challenge, computational cognitive scientists have developed models of *joint intentionality* [73], where cooperating agents conceive of themselves as a *group agent* with a common goal in mind [66, 69]. This approach allows groups of agents to converge on shared goals [70] and efficient task decompositions [78] in a decentralized manner with limited recursion. CLIPS uses this insight to model joint planning and pragmatic instruction generation in humans.

Value Alignment and Assistance Games. CLIPS is a solution strategy for language-augmented goal assistance games, an extension of the assistance game formalism for human-AI value alignment [25, 58]. In contrast to prior work, our framework leverages joint intentionality when modeling human principals, thereby requiring less recursion than iterated best-response [25], while avoiding the intractability of equilibrium solutions [18, 49].

Instruction Following with Language Models. CLIPS is a form of grounded instruction following from natural language [50, 65, 71], using LLMs to score the likelihood of an utterance given a grounded command. While many recent studies have employed LLMs for translating language into actions [2, 64] or task specifications [41, 46, 82], we opt for an explicitly Bayesian approach, allowing our method to integrate information from both actions and instructions in a principled, reliable, and modular fashion.

5 DISCUSSION AND FUTURE WORK

In this paper, we introduced cooperative language-guided inverse plan search (CLIPS) as a Bayesian architecture for pragmatic instruction following and goal assistance. By using a structured cognitive model of how speakers produce both actions and utterances given their goals, CLIPS is able to integrate information from observed actions and ambiguous instruction, inferring a distribution over speaker’s goals and intentions. Using this distribution, a CLIPS assistant is then able to assist the speaker in achieving their goal through expected cost minimization. Through our experiments, we show that the proposed architecture produces human-like outputs on goal inference and assistance tasks. Compared with the baselines, CLIPS also makes more accurate goal inferences, and chooses actions that result in more efficient plans.

While CLIPS demonstrates compelling theoretical properties and strong empirical performance, a number of challenges remain before it can be applied to scenarios with a larger number of goals and assistance options. As noted in Section 2.2, the Bayesian inference strategy we adopt is fully enumerative, but this approach breaks down once the space of possible goals, plans, and commands grows sufficiently large. Recent advances in probabilistic programming could overcome these bottlenecks. For example, more sophisticated sequential Monte Carlo strategies could be used to flexibly propose and update hypotheses, focusing computation on only the most likely sets of latent variables [42, 86]. These strategies can be applied to perform *constrained decoding* from large language models [43, 76], thereby allowing LLMs to be used as sound proposal distributions over a grammar for commands. This would benefit from LLMs’ performance at few-shot semantic parsing and translation [62], while preserving the Bayesian nature of our framework.

The assistance policy we introduced in Algorithm 2 could also be extended in a number of ways. In many assistive settings, uncertainty is sufficiently high that it makes sense for the assistant to take *information gathering actions*. Indeed, this is what people typically do when we are confused: we ask questions. To enable this ability, assistants should perform *belief-space planning* over the set of possible goal beliefs [47, 68], thinking ahead about which actions reveal more information about the principal’s goals. Such planning would alleviate issues that can arise due to the Q_{MDP} approximation [27, 45] used by Algorithm 2. In the longer run, assistants could be augmented with natural language outputs, enabling them to ask clarifying questions by planning ahead in a white-box, interpretable manner. If successful, this would constitute a considerable step towards trustworthy assistive AI that effectively collaborates and communicates with humans.

ACKNOWLEDGMENTS

This work was funded in part by the DARPA Machine Common Sense, AFOSR, and ONR Science of AI programs, along with the MIT-IBM Watson AI Lab, and gifts from Reid Hoffman and the Siegel Family Foundation. Tan Zhi-Xuan is supported by an Open Philanthropy AI Fellowship.

Errata: Along with typographical fixes, this version corrects minor errors in the GPT-4V assistance results that are present in the version of this paper published in the ACM Digital Library.

REFERENCES

- [1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. GPT-4 Technical Report. *arXiv preprint arXiv:2303.08774* (2023).
- [2] Michael Ahn, Anthony Brohan, Noah Brown, Yevgen Chebotar, Omar Cortes, Byron David, Chelsea Finn, Chuyuan Fu, Keerthana Gopalakrishnan, Karol Hausman, et al. 2022. Do as I can, not as I say: Grounding language in robotic affordances. *arXiv preprint arXiv:2204.01691* (2022).
- [3] Arwa Alanqary, Gloria Z Lin, Joie Le, Tan Zhi-Xuan, Vikash K Mansinghka, and Joshua B Tenenbaum. 2021. Modeling the mistakes of boundedly rational agents within a Bayesian theory of mind. *arXiv preprint arXiv:2106.13249* (2021).
- [4] John Langshaw Austin. 1975. *How to do things with words*. Vol. 88. Oxford University Press.
- [5] Davide Aversa, Sebastian Sardina, and Stavros Vassos. 2016. Pruning and preprocessing methods for inventory-aware pathfinding. In *2016 IEEE Conference on Computational Intelligence and Games (CIG)*. IEEE, 1–8.
- [6] Chris Baker, Rebecca Saxe, and Joshua Tenenbaum. 2011. Bayesian Theory of Mind: Modeling Joint Belief-Desire Attribution. In *Proceedings of the Annual Meeting of the Cognitive Science Society*, 33 (33).
- [7] Chris L Baker, Julian Jara-Ettinger, Rebecca Saxe, and Joshua B Tenenbaum. 2017. Rational quantitative attribution of beliefs, desires and percepts in human mentalizing. *Nature Human Behaviour* 1, 4 (2017), 1–10.
- [8] Chris L Baker, Rebecca Saxe, and Joshua B Tenenbaum. 2009. Action understanding as inverse planning. *Cognition* 113, 3 (2009), 329–349.
- [9] Andrew G Barto, Steven J Bradtke, and Satinder P Singh. 1995. Learning to act using real-time dynamic programming. *Artificial intelligence* 72, 1-2 (1995), 81–138.
- [10] Craig Boutilier. 1996. Planning, learning and coordination in multiagent decision processes. In *TARK*, Vol. 96. Citeseer, 195–210.
- [11] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems* 33 (2020), 1877–1901.
- [12] Luca M. Schulze Buschoff, Elif Akata, Matthias Bethge, and Eric Schulz. 2024. Visual cognition in multimodal large language models. *arXiv:2311.16093 [cs.LG]*
- [13] Colin F Camerer, Teck-Hua Ho, and Juin-Kuan Chong. 2004. A cognitive hierarchy model of games. *The Quarterly Journal of Economics* 119, 3 (2004), 861–898.
- [14] Jaedug Choi and Kee-Eung Kim. 2014. Hierarchical bayesian inverse reinforcement learning. *IEEE transactions on cybernetics* 45, 4 (2014), 793–805.
- [15] Marco F Cusumano-Towner, Feras A Saad, Alexander K Lew, and Vikash K Mansinghka. 2019. Gen: a general-purpose probabilistic programming system with programmable inference. In *Proceedings of the 40th ACM SIGPLAN Conference on Programming Language Design and Implementation*. ACM, 221–236.
- [16] Pierre Del Moral, Arnaud Doucet, and Ajay Jasra. 2006. Sequential Monte Carlo samplers. *Journal of the Royal Statistical Society Series B: Statistical Methodology* 68, 3 (2006), 411–436.
- [17] Prashant Doshi and Piotr J Gmytrasiewicz. 2009. Monte Carlo sampling methods for approximating interactive POMDPs. *Journal of Artificial Intelligence Research* 34 (2009), 297–337.
- [18] Jaime F Fisac, Monica A Gates, Jessica B Hamrick, Chang Liu, Dylan Hadfield-Menell, Malayandi Palaniappan, Dhruv Malik, S Shankar Sastry, Thomas L Griffiths, and Anca D Dragan. 2017. Pragmatic-Pedagogic Value Alignment. *arXiv preprint arXiv:1707.06354* (2017).
- [19] Daniel Fried, Jacob Andreas, and Dan Klein. 2018. Unified Pragmatic Models for Generating and Following Instructions. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, Marilyn Walker, Heng Ji, and Amanda Stent (Eds.). Association for Computational Linguistics, New Orleans, Louisiana, 1951–1963. <https://doi.org/10.18653/v1/N18-1177>
- [20] Justin Fu, Anoop Korattikara, Sergey Levine, and Sergio Guadarrama. 2019. From language to goals: Inverse reinforcement learning for vision-based instruction following. *arXiv preprint arXiv:1902.07742* (2019).
- [21] György Gergely and Gergely Csibra. 2003. Teleological reasoning in infancy: The naive theory of rational action. *Trends in cognitive sciences* 7, 7 (2003), 287–292.
- [22] Rachel Giora. 2004. On the graded salience hypothesis. (2004).
- [23] Noah D Goodman and Michael C Frank. 2016. Pragmatic language interpretation as probabilistic inference. *Trends in cognitive sciences* 20, 11 (2016), 818–829.
- [24] Herbert P Grice. 1975. Logic and Conversation. In *Speech Acts*. Brill, 41–58.
- [25] Dylan Hadfield-Menell, Stuart J Russell, Pieter Abbeel, and Anca Dragan. 2016. Cooperative inverse reinforcement learning. In *Advances in neural information processing systems*. 3909–3917.
- [26] Patrik Haslum, Blai Bonet, Héctor Geffner, et al. 2005. New admissible heuristics for domain-independent planning. In *AAAI*, Vol. 5. 9–13.
- [27] Milos Hauskrecht. 2000. Value-function approximations for partially observable Markov decision processes. *Journal of artificial intelligence research* 13 (2000), 33–94.
- [28] Michael Held and Richard M Karp. 1971. The traveling-salesman problem and minimum spanning trees: Part II. *Mathematical programming* 1, 1 (1971), 6–25.
- [29] Malte Helmert and Carmel Domshlak. 2009. Landmarks, critical paths and abstractions: what’s the difference anyway?. In *Proceedings of the International Conference on Automated Planning and Scheduling*, Vol. 19. 162–169.
- [30] Carlos Hernández, Tansel Uras, Sven Koenig, Jorge A Baier, Xiaoxun Sun, and Pedro Meseguer. 2015. Reusing cost-minimal paths for goal-directed navigation in partially known terrains. *Autonomous Agents and Multi-Agent Systems* 29 (2015), 850–895.
- [31] Aspen K Hopkins, Alex Renda, and Michael Carbin. 2023. Can LLMs Generate Random Numbers? Evaluating LLM Sampling in Controlled Domains. In *ICML 2023 Workshop: Sampling and Optimization in Discrete Space*.
- [32] Julian Jara-Ettinger. 2019. Theory of mind as inverse reinforcement learning. *Current Opinion in Behavioral Sciences* 29 (2019), 105–110.
- [33] Julian Jara-Ettinger, Laura Schulz, and Josh Tenenbaum. 2019. The Naive Utility Calculus as a unified, quantitative framework for action understanding. *PsyArXiv* (2019).
- [34] Hong Jun Jeon, Smitha Milli, and Anca Dragan. 2020. Reward-rational (implicit) choice: A unifying formalism for reward learning. *Advances in Neural Information Processing Systems* 33 (2020), 4415–4426.
- [35] Chuanyang Jin, Yutong Wu, Jing Cao, Jiannan Xiang, Yen-Ling Kuo, Zhiting Hu, Tomer Ullman, Antonio Torralba, Joshua B Tenenbaum, and Tianmin Shu. 2024. MMTom-QA: Multimodal Theory of Mind Question Answering. *arXiv preprint arXiv:2401.08743* (2024).
- [36] Istvan Kecskes et al. 2004. The role of salience in processing pragmatic units. *Acta Linguistica Hungarica (Since 2017 Acta Linguistica Academica)* 51, 3-4 (2004), 309–324.
- [37] Robert Kirk, Ishita Mediratta, Christoforos Nalmpantis, Jelena Luketina, Eric Hambro, Edward Grefenstette, and Roberta Raileanu. 2023. Understanding the effects of RLHF on LLM generalisation and diversity. *arXiv preprint arXiv:2310.06452* (2023).
- [38] Sven Koenig and Maxim Likhachev. 2006. Real-Time Adaptive A*. In *Proceedings of the fifth international joint conference on Autonomous agents and multiagent systems*. 281–288.
- [39] Sven Koenig and Xiaoxun Sun. 2009. Comparing real-time and incremental heuristic search for real-time situated agents. *Autonomous Agents and Multi-Agent Systems* 18 (2009), 313–341.
- [40] Richard E Korf. 1990. Real-time heuristic search. *Artificial intelligence* 42, 2-3 (1990), 189–211.
- [41] Minae Kwon, Sang Michael Xie, Kalesha Bullard, and Dorsa Sadigh. 2023. Reward Design with Language Models. In *The Eleventh International Conference on Learning Representations*. <https://openreview.net/forum?id=10uNUg15Kl>
- [42] Alexander K Lew, George Matheos, Tan Zhi-Xuan, Matin Ghavamizadeh, Nishad Gothoskar, Stuart Russell, and Vikash K Mansinghka. 2023. SMCPC3: Sequential Monte Carlo with probabilistic program proposals. In *International Conference on Artificial Intelligence and Statistics*. PMLR, 7061–7088.
- [43] Alexander K Lew, Tan Zhi-Xuan, Gabriel Grand, and Vikash K Mansinghka. 2023. Sequential Monte Carlo Steering of Large Language Models using Probabilistic Programs. *arXiv preprint arXiv:2306.03081* (2023).
- [44] Jessy Lin, Daniel Fried, Dan Klein, and Anca Dragan. 2022. Inferring rewards from language in context. *arXiv preprint arXiv:2204.02515* (2022).
- [45] Michael L Littman, Anthony R Cassandra, and Leslie Pack Kaelbling. 1995. Learning policies for partially observable environments: Scaling up. In *Machine Learning Proceedings 1995*. Elsevier, 362–370.
- [46] Jason Xinyu Liu, Ziyi Yang, Benjamin Schornstein, Sam Liang, Ifrah Idrees, Stefanie Tellex, and Ankit Shah. 2022. Lang2LTL: Translating Natural Language Commands to Temporal Specification with Large Language Models. In *Workshop on Language and Robotics at CoRL 2022*.
- [47] Owen Macindoe, Leslie Pack Kaelbling, and Tomás Lozano-Pérez. 2012. POMCoP: Belief Space Planning for Sidekicks in Cooperative Games. In *AAAI*.
- [48] Drew McDermott, Malik Ghallab, Adele Howe, Craig Knoblock, Ashwin Ram, Manuela Veloso, Daniel Weld, and David Wilkins. 1998. PDDL - The Planning Domain Definition Language.
- [49] Smitha Milli and Anca D Dragan. 2020. Literal or pedagogic human? analyzing human model misspecification in objective learning. In *Uncertainty in artificial intelligence*. PMLR, 925–934.
- [50] Dipendra K Misra, Jaeyong Sung, Kevin Lee, and Ashutosh Saxena. 2016. Tell me Dave: Context-sensitive grounding of natural language to manipulation instructions. *The International Journal of Robotics Research* 35, 1-3 (2016), 281–300.
- [51] Kevin Murphy and Stuart Russell. 2001. Rao-Blackwellized particle filtering for dynamic Bayesian networks. In *Sequential Monte Carlo methods in practice*. Springer, 499–515.
- [52] Truong-Huy Nguyen, David Hsu, Wee-Sun Lee, Tze-Yun Leong, Leslie Kaelbling, Tomas Lozano-Perez, and Andrew Grant. 2011. Capir: Collaborative action planning with intention recognition. In *Proceedings of the AAAI Conference on Artificial Intelligence and Interactive Digital Entertainment*, Vol. 7. 61–66.

- [53] Judea Pearl. 2012. The do-calculus revisited. In *Proceedings of the Twenty-Eighth Conference on Uncertainty in Artificial Intelligence*. 3–11.
- [54] Xavier Puig, Kevin Ra, Marko Boben, Jiaman Li, Tingwu Wang, Sanja Fidler, and Antonio Torralba. 2018. Virtualhome: Simulating household activities via programs. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 8494–8502.
- [55] Xavier Puig, Tianmin Shu, Shuang Li, Zilin Wang, Yuan-Hong Liao, Joshua B Tenenbaum, Sanja Fidler, and Antonio Torralba. 2020. Watch-and-help: A challenge for social perception and human-ai collaboration. *arXiv preprint arXiv:2010.09890* (2020).
- [56] Deepak Ramachandran and Eyal Amir. 2007. Bayesian Inverse Reinforcement Learning. In *IJCAI*, Vol. 7. 2586–2591.
- [57] Miguel Ramírez and Hector Geffner. 2010. Probabilistic plan recognition using off-the-shelf classical planners. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 24.
- [58] Stuart Russell. 2021. Human-Compatible Artificial Intelligence. *Human-like machine intelligence* (2021), 3–23.
- [59] Sandhya Saisubramanian, Kyle Hollins Wray, Luis Pineda, and Shlomo Zilberstein. 2019. Planning in stochastic environments with goal uncertainty. In *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 1649–1654.
- [60] Melanie Sclar, Sachin Kumar, Peter West, Alane Suhr, Yejin Choi, and Yulia Tsvetkov. 2023. Minding Language Models’ (Lack of) Theory of Mind: A Plug-and-Play Multi-Character Belief Tracker. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (Eds.). Association for Computational Linguistics, Toronto, Canada, 13960–13980. <https://doi.org/10.18653/v1/2023.acl-long.780>
- [61] Ankit Shah, Shen Li, and Julie Shah. 2020. Planning with uncertain specifications (puns). *IEEE Robotics and Automation Letters* 5, 2 (2020), 3414–3421.
- [62] Richard Shin, Christopher H Lin, Sam Thomson, Charles Chen, Subhro Roy, Emmanouil Antonios Platanios, Adam Pauls, Dan Klein, Jason Eisner, and Benjamin Van Durme. 2021. Constrained language models yield few-shot semantic parsers. *arXiv preprint arXiv:2104.08768* (2021).
- [63] Michael Shum, Max Kleiman-Weiner, Michael L Littman, and Joshua B Tenenbaum. 2019. Theory of minds: Understanding behavior in groups through inverse planning. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 33. 6163–6170.
- [64] Ishika Singh, Valts Blukis, Arsalan Mousavian, Ankit Goyal, Danfei Xu, Jonathan Tremblay, Dieter Fox, Jesse Thomason, and Animesh Garg. 2023. Progprompt: Generating situated robot task plans using large language models. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 11523–11530.
- [65] Shawn Squire, Stefanie Tellex, Dilip Arumugam, and Lei Yang. 2015. Grounding English commands to reward functions. In *Robotics: Science and Systems*.
- [66] Stephanie Stacy, Chenfei Li, Minglu Zhao, Yiling Yun, Qingyi Zhao, Max Kleiman-Weiner, and Tao Gao. 2021. Modeling communication to coordinate perspectives in cooperation. *arXiv preprint arXiv:2106.02164* (2021).
- [67] Theodore R Sumers, Mark K Ho, Thomas L Griffiths, and Robert D Hawkins. 2023. Reconciling truthfulness and relevance as epistemic and decision-theoretic utility. *Psychological Review* (2023).
- [68] Zachary Sunberg and Mykel Kochenderfer. 2018. Online algorithms for POMDPs with continuous state, action, and observation spaces. In *Proceedings of the International Conference on Automated Planning and Scheduling*, Vol. 28. 259–263.
- [69] Ning Tang, Siyi Gong, Minglu Zhao, Chenya Gu, Jifan Zhou, Mowei Shen, and Tao Gao. 2022. Exploring an Imagined “We” in human collective hunting: Joint commitment within shared intentionality. In *Proceedings of the annual meeting of the cognitive science society*, Vol. 44.
- [70] Ning Tang, Stephanie Stacy, Minglu Zhao, Gabriel Marquez, and Tao Gao. 2020. Bootstrapping an Imagined We for Cooperation. In *CogSci*.
- [71] Stefanie Tellex, Thomas Kollar, Steven Dickerson, Matthew Walter, Ashis Banerjee, Seth Teller, and Nicholas Roy. 2011. Understanding natural language commands for robotic navigation and mobile manipulation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 25. 1507–1514.
- [72] William R Thompson. 1935. On the theory of apportionment. *American Journal of Mathematics* 57, 2 (1935), 450–456.
- [73] Michael Tomasello and Malinda Carpenter. 2007. Shared intentionality. *Developmental science* 10, 1 (2007), 121–125.
- [74] Hsiao-Yu Tung, Adam W Harley, Liang-Kang Huang, and Katerina Fragkiadaki. 2018. Reward learning from narrated demonstrations. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 7004–7013.
- [75] Yanming Wan, Jiayuan Mao, and Josh Tenenbaum. 2022. HandMeThat: Human-Robot Communication in Physical and Social Environments. *Advances in Neural Information Processing Systems* 35 (2022), 12014–12026.
- [76] Brandon T Willard and Rémi Louf. 2023. Efficient Guided Generation for LLMs. *arXiv preprint arXiv:2307.09702* (2023).
- [77] Edward C Williams, Nakul Gopalan, Mine Rhee, and Stefanie Tellex. 2018. Learning to parse natural language to grounded reward functions with weak supervision. In *2018 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 4430–4436.
- [78] Sarah A Wu, Rose E Wang, James A Evans, Joshua B Tenenbaum, David C Parkes, and Max Kleiman-Weiner. 2021. Too Many Cooks: Bayesian Inference for Coordinating Multi-Agent Collaboration. *Topics in Cognitive Science* 13, 2 (2021), 414–432.
- [79] Frank Yates. 1948. Systematic sampling. *Philosophical Transactions of the Royal Society of London. Series A, Mathematical and Physical Sciences* 241, 834 (1948), 345–377.
- [80] Lance Ying, Katherine M Collins, Megan Wei, Cedegao E Zhang, Tan Zhi-Xuan, Adrian Weller, Joshua B Tenenbaum, and Lionel Wong. 2023. The Neuro-Symbolic Inverse Planning Engine (NIPE): Modeling Probabilistic Social Inferences from Linguistic Inputs. *arXiv preprint arXiv:2306.14325* (2023).
- [81] Lance Ying, Tan Zhi-Xuan, Vikash Mansinghka, and Joshua B Tenenbaum. 2023. Inferring the goals of communicating agents from actions and instructions. In *Proceedings of the AAAI Symposium Series*, Vol. 2. 26–33.
- [82] Wenhao Yu, Nimrod Gileadi, Chuyuan Fu, Sean Kirmani, Kuang-Huei Lee, Montse Gonzalez Arenas, Hao-Tien Lewis Chiang, Tom Erez, Leonard Hasenclever, Jan Humplik, et al. 2023. Language to Rewards for Robotic Skill Synthesis. *arXiv preprint arXiv:2306.08647* (2023).
- [83] Tan Zhi-Xuan. 2022. GenGPT3.jl: GPT-3 as a generative function in Gen.jl. <https://github.com/probcomp/GenGPT3.jl>
- [84] Tan Zhi-Xuan. 2022. *PDDL.jl: An Extensible Interpreter and Compiler Interface for Fast and Flexible AI Planning*. Ph.D. Dissertation. Massachusetts Institute of Technology.
- [85] Tan Zhi-Xuan, Nishad Gothoskar, Falk Pollok, Dan Gutfreund, Joshua B Tenenbaum, and Vikash K Mansinghka. 2022. Solving the Baby Intuition Benchmark with a hierarchically Bayesian theory-of-mind. *arXiv preprint arXiv:2208.02914* (2022).
- [86] Tan Zhi-Xuan, Jordyn Mann, Tom Silver, Josh Tenenbaum, and Vikash Mansinghka. 2020. Online Bayesian Goal Inference for Boundedly Rational Planning Agents. *Advances in Neural Information Processing Systems* 33 (2020).
- [87] Brian D Ziebart, Andrew L Maas, J Andrew Bagnell, and Anind K Dey. 2008. Maximum entropy inverse reinforcement learning. In *AAAI*, Vol. 8. Chicago, IL, USA, 1433–1438.

A GOAL ASSISTANCE PROBLEMS

Table A1 and Table A2 describe the scenarios and conditions for our goal assistance problems in multi-agent Doors, Keys & Gems (mDKG) and VirtualHome (VH) respectively. VH scenarios were encoded in PDDL [48] for use with our inverse planning system, but they can be easily be translated to VirtualHome’s evolving graph simulator or Unity simulator [54].

Table A1: Goal assistance problems in Doors, Keys & Gems. 30 problems were divided into those that required assistance with either picking up keys or unlocking doors. Actions and utterances were varied according to the following factors.

Factor	Condition	Description
Instruction Uncertainty	Ambiguous	Ambiguous instruction that can be associated with many possible items. (e.g. “Can you get me the red key?” when multiple red keys are available.)
	Partial	The principal only requests a subset of the items they require to achieve their goal. (e.g. “Can you get me the red key?” when both red and blue keys to reach its goal.)
Agent Involvement	Assistant Actions	The instruction involves only actions that the assistant should take. (e.g. “Can you pick up that key?”, “Please unlock the red doors.”)
	Joint Actions	The instructions contain information about both agents’ actions. (e.g. “I’m getting the red key, can you get the blue key?”)
Goal Uncertainty	Low Uncertainty	The actions and instructions combined are largely sufficient to uniquely determine the goal. (e.g. In Figure 4b, only the red gem is likely.)
	High Uncertainty	The principal’s goal remains ambiguous even after observing both their actions and instructions. (e.g. In Figure 4d, the three gems on the left are plausible goals.)

Table A2: Goal assistance problems in VirtualHome. 25 problems were divided into (a) three scenarios, which correspond to different initial states of the environment, and had different sets of possible goals. Actions and utterances (b) were varied across several conditions.

Scenario	Goal Category (k)	Description
Food Preparation (6 goals)	Salads (3)	Place a chef’s knife & ingredients (onion, cucumber/potato & tomato/chicken) on the table.
	Stews (3)	Place an onion, carrot, a main (salmon, chicken, or potato) & wine on the table.
Table Setting (12 goals)	Plates (4)	Place n plates on the table, where $n \in [1, 4]$.
	Plates & Cutlery (4)	Place n plates, n forks & n knives on the table, where $n \in [1, 4]$.
	Plates & Bowls (4)	Place n plates, n forks, n knives & n bowls on the table, where $n \in [1, 4]$.
Drink Pairing (12 goals)	Wine (3)	Place wine and n wineglasses on the table, $n \in [1, 3]$.
	Wine Pairing (3)	Place wine, cheese, n wineglasses & n forks on the table, $n \in [1, 3]$.
	Juice (3)	Place juice & n waterglasses on the table, $n \in [1, 3]$.
	Juice Pairing (3)	Place juice, n waterglasses, n cupcakes & n forks on the table, $n \in [1, 3]$.

(a) Scenarios and goals.

Factor	Condition	Description
Instruction Uncertainty	Ambiguous Number	Instruction which mentions an ambiguous number of items to be collected. (e.g. “Can you help with these plates?”)
	Ambiguous Type	Instruction which is ambiguous about the type of item to be collected. (e.g. “Can you get two glasses?”, “And I need the veggies for the stew.”)
Utterance Count	Single	The principal gives a single instruction to the assistant. (e.g. “Go get the forks and knives.”, “Can you get some wine?”)
	Multiple	The principal gives an instruction to the assistant, takes a few more actions, then gives another instruction. (e.g. “Could you find two glasses?” followed by “Get some forks as well.”)

(b) Factors and conditions.

B MODEL CONFIGURATION

In this section, we provide further details about how CLIPS and each of our baselines was configured for each domain.

B.1 CLIPS & Unimodal Inverse Planning

We describe the configuration for CLIPS and unimodal inverse planning together, since the only difference between the former and the latter is whether both action likelihoods $P(a_t^h | s_t, \pi)$ and utterance likelihoods $P(u_t | s_t, \pi)$ are incorporated into the weight w associated with each hypothesized goal g , cost function C and policy π in Algorithm 1.

Goal and Cost Priors. In both domains, we assumed a uniform goal prior $P(g)$ over the set of possible goals (4 in mDKG, 6–12 in VH). In mDKG, we considered four possible cost functions $C \in \mathcal{C}$, and placed a uniform prior $P(C)$ over such cost functions. This made CLIPS more robust to the possibility that the human would prefer to do less (or more) work, asking the robot to pick up keys or unlock doors that would be sub-optimal if action costs were equal between both agents. The costs for each action were as follows:

- Human pick-up: 5. Robot pick-up: 1. Wait cost: 0.6. All other action costs: 1.
- Human pick-up: 5. Robot pick-up: 1. Wait cost: 0.6. Other human costs: 2. Other robot costs: 1.
- Human pick-up: 5. Robot pick-up: 1. Robot unlock: 5. Wait cost: 0.6. All other action costs: 1.
- Human pick-up: 5. Robot pick-up: 1. Robot unlock: 5. Wait cost: 0.6. Other human costs: 2. Other robot costs: 1.

In VH, we avoided the need to perform online inference over multiple cost functions by using a slightly more flexible utterance model that allowed for the possibility that the principal might ask the assistant to carry out sub-optimal actions (e.g. asking the assistant to acquire *all* remaining items, even though the optimal joint plan would have both agents share the work more equally). This allowed us to use a *single* cost function that assigned unit cost to all actions for both agents. Future work should investigate how more flexible models for multi-agent task division (e.g. [78]) could be integrated with CLIPS, along with offline learning of cost function priors.

Policy Configuration. In Bayesian inverse planning, an algorithm is required to compute the joint policy π for each possible goal. While it is possible to do this using a domain-general planning algorithm and search heuristic, we configured our planners in a domain-specific manner to achieve greater speed and reliability.

Doors, Keys & Gems. In mDKG, we used the Real-Time Adaptive A* (RTAA*) variant [38] of real-time heuristic search [40] to incrementally compute the value function $V(s)$ for the policy π . Given $V(s)$, one can then compute $Q(s, a) = C(s, a) + \sum_{s'} P(s' | s, a) V(s')$. Each time `RTHS-POLICY-UPDATE` is called, we perform A* search up to a search budget of 2^{18} nodes from each state s_n that neighbors the current state s , using the current estimate \hat{V} of the value function as the search heuristic. Once A* search terminates at state s_f (either by reaching the goal or using up the budget), we update the value function for *all* states s_i in the interior of the search tree as follows:

$$\hat{V}(s_i) \leftarrow \hat{V}(s_f) - C(s_n, s_f) + C(s_n, s_i)$$

where $C(s_n, s_i)$ is the cost of getting from s_n to s_i by traversing the search tree. As shown in [38], as long as \hat{V} is initialized to a *consistent* (under-)estimate of the true value function, subsequent \hat{V} values will remain consistent. We ensured this by (lazily) initializing \hat{V} with a minimum spanning tree heuristic [28] over the graph of keys and doors required to reach the goal, using a preprocessing step used in inventory-aware pathfinding [5] to extract necessary keys and doors. To further accelerate planning, we pruned actions that were irrelevant to the goal (e.g. picking up unnecessary keys), and implemented cost-minimal path reuse from the tree-adaptive variant of RTAA* [30], allowing us to terminate A* search whenever an existing path to the goal was found. Collectively, these algorithmic choices allowed us to estimate the optimal value function V^* in close to real-time for most timesteps, despite the difficulty of optimal multi-agent planning.

VirtualHome. In VH, we also experimented with using RTAA* to compute the policy π . However, this unpredictably increased the latency of inverse planning when the search budget was too high, and led to unstable Q -value estimates when the search budget was too low. Instead we found that it was sufficient to compute an analytic lower bound \hat{V} to the optimal value function V^* by counting the number of remaining items and movements required to achieve the goal. Using \hat{V} was both faster (no search was required), and produced more consistent estimates of the relative Q -values $Q(s, a) - V(s)$ (also known as the *advantage function*) for each state s and goal g , resulting in more stable action likelihoods $P(a | s, \pi)$ and goal inferences. However, this required hand-deriving a lower bound on V^* . In the future, we hope to explore the use of domain-general lower bounds on V^* using admissible planning heuristics [26, 29], and to implement more consistent value function updates based on the LSS-LRTA* variant of real-time heuristic search [39], allowing us to stably estimate V^* with a low search budget.

Rationality Modeling. As in hierarchical Bayesian approaches to inverse reinforcement learning [14] and theory-of-mind [85], we placed a Gamma prior over the rationality parameter β for the Boltzmann policy, discretized across 33 values from 2^{-3} to 2^5 by incrementing the exponent with a step size of 0.25. In mDKG, we used a $\text{Gamma}(0.5, 1)$ prior, while in VH, we used a $\text{Gamma}(1, 2)$ to match the less optimal actions taken in that domain. For efficient inference, we used Rao-Blackwellization [51], marginalizing over all values of β when computing the action likelihood $P(a^h | s_t, \pi)$, then analytically updating the local posterior $P(\beta | a_{1:t}^h, s_{1:t}, \pi, g)$.

Utterance Model. In the CLIPS utterance model, salient actions $\alpha_{1:m}$ are first extracted from the joint policy π , before being sub-sampled into a command c , which is then translated into an utterance u . In our implementation, we enumerated over commands to reduce variance. This section describes how we extracted salient actions, the structure of commands, and the command enumeration strategy for each domain.

Extracting Salient Actions. Following work in linguistics [22, 36] and RSA theory [23], we assumed that some types of actions were more likely to be communicated because of their *salience* to both the speaker and the listener. To model this, after rolling out the joint policy π for up to H steps, we filtered out all non-salient actions, leaving only the salient actions $\alpha_{1:m}$ to be communicated. In VH, to account for the possibility of the principal delegating *all* remaining items to the assistant (instead of equally sharing the load), we enforced that the principal could not collect new items when rolling out the policy.

As in [81], we manually defined action salience: In mDKG, pickup, unlock, and handover actions and iscolor predicates were salient, while in VH, grab actions and itemtype predicates were salient. In general, however, action salience could be learned, either from a corpus, or by offline estimation of which actions tend to be more *unique* and hence *informative* for a particular goal and plan.

Command Enumeration. Formally, we define a command $c := ([a, \dots], \{p, \dots\}, \{(v, \tau), \dots\})$ as a sequence of (possibly lifted) actions $[a, \dots]$, a set of associated predicates $\{p, \dots\}$ that qualify an action’s arguments, and a set $\{(v, \tau), \dots\}$ of free variables v with associated types τ which appear as an argument in at least one action in $[a, \dots]$. For example, the command (handover robot human ?key1) where (iscolor ?key1 blue) has one partially lifted action $a = (\text{handover robot human ?key1})$, along with a predicate $p = (\text{iscolor ?key1 blue})$ that describes the variable $v = \text{?key1}$ as blue, where the type of ?key1 is $\tau = \text{key}$.

To generate commands from a set of salient actions $\alpha_{1:m}$ and associated predicates $p_{1:m}$, we enumerated over all subsets of salient action-predicate pairs of up to size $K = 3$ in mDKG and $K = 4$ in VH, before pruning away combinations of actions that were less likely. In mDKG, we did this by excluding commands with more than two distinct action types, commands where *all* actions were taken by the speaker/principal, and also commands which involved chains of dependent actions (e.g. picking up a key, followed by handing over the same key). In VH, we excluded all commands that involved *any* action by the speaker (since our VH dataset did not contain joint instructions). Following enumeration of these (ground) commands, we then lifted the commands, replacing concrete object names with corresponding typed variables. We also replaced all agent identifiers such as human and robot with corresponding indexicals (me and you). Duplicate lifted commands were then removed.

These choices were made partly to limit the number of commands enumerated. Future extensions to CLIPS that make use of constrained LLM decoding [43, 62] could avoid the limitations involved in command enumeration, directly translating an utterance u into a much larger space of possible commands (e.g. a context-free grammar over commands).

Language Models. In mDKG, we used OpenAI’s 6.7B curie model [11] as our LLM⁵, while in VH, we used their more recently released davinci-002 model. Both of these were base models, providing more diffuse probabilities that avoid issues with over-confidence and mode collapse that can arise in LLMs finetuned on instructions or human preferences [31, 37].

LLMs were instantiated as generative functions in the Gen probabilistic programming system [15] using the GenGPT3.jl library [83], allowing us to easily score the probability $P(u|c, \mathcal{E})$ of an utterance u given a command c and few-shot examples \mathcal{E} provided in the prompt. Each utterance was scored against all enumerated commands in parallel via a batched API request. In mDKG, we provided 18 examples of commands being translated to utterances, while in VH, we provided 12 examples. Examples can be found in the source code.

Assistance Policy. When running the Q_{MDP} assistance policy shown in Algorithm 2, we filtered out goal hypotheses with a weight w^i of less than 0.02, so as to avoid the cost of updating the policies π^i for unlikely goals. In mDKG, we updated the policies using RTAA* with a node budget 2^{16} and a time limit of 10s. In VH, we updated Q -values using the analytic lower bound \hat{V} mentioned earlier.

We evaluated the assistance policy against a simulated human principal, which was designed to (i) follow a “ground truth” plan whenever possible, (ii) switch to following an approximately optimal joint policy when the following the original plan was not possible, (iii) fallback to a single-agent policy upon detection that the assistant was not following the joint policy to the true goal (i.e. not being helpful). To detect this, the principal compared how likely the assistant was to be following a joint policy to the true goal (Boltzmann, with $\beta = 1$) vs. a random baseline, switching to the fallback policy once it was 10 times more likely for the assistant to be acting at random.

B.2 LLM-Based Literal Listener

Our literal listener baselines consisted of two components: literal command inference, and literal assistance. We describe each below.

Literal Command Inference. Command inference was implemented similarly to inference over the CLIPS utterance model. We first enumerated over a set of possible commands in the current state s , then used a LLM to score the probability $P(u|c, s)$ of the observed utterance u for each command c . Normalizing these probabilities produced a posterior over commands $P(c|u, s)$.

Utterances. In mDKG, we made command inference easier by replacing joint instructions (e.g. “I’m getting the red key, can you get the blue key?”) with assistant-directed instructions (e.g. “Can you get the blue key?”). In VH, we handled multiple utterances u_{t_1}, u_{t_2} by merging the inferred commands $c = c_{t_1} \cup c_{t_2}$. This created a product distribution $P(c|u_{t_1}, u_{t_2}, s_{t_1}, s_{t_2}) = P(c_{t_1}|u_{t_1}, s_{t_1})P(c_{t_2}|u_{t_2}, s_{t_2})$ over merged commands.

Command Enumeration. To enumerate commands, we first enumerated all salient actions that were reachable from the current state s . We then enumerated all subsets the set of salient actions of up to size $K = 3$ in mDKG and $K = 4$, pruning away certain combinations as in the CLIPS utterance model. Commands were then lifted, and duplicate commands were removed.

⁵Unfortunately, as of 4 January 2024, OpenAI has shutdown access to curie and the original GPT-3 family of models, making it no longer possible to reproduce our results exactly. Nonetheless, we expect similar results when using davinci-002 as a replacement model.

Language Models. We used `text-davinci-002` for mDKG. Given the larger space of commands for the literal listener, we found that using a RLHF finetuned model such as `text-davinci-002` ensured that more posterior probability mass concentrated on the intuitively correct set of commands. To simulate this effect in our VH experiments after `text-davinci-002` access was shutdown by OpenAI, we used `davinci-002` with a temperature of 0.5. 12 few-shot examples were used for both mDKG and VH.

Literal Assistance. Literal assistance involved several steps: (i) Systematically sampling a command c from the inferred distribution $P(c|u, s)$; (ii) Converting the command c to an equivalent goal formula g_c ; (iii) Generating an assistive plan that first achieved g_c , then the true goal g .

Sampling Commands. As described in the main text, we estimated the *mean* performance of literal assistance by drawing $m = 10$ samples from the distribution $P(c|u, s)$ and averaging our metrics across these samples. To do this in a low-variance manner, we first sorted the commands by their inferred probabilities. We then used systematic sampling [79], selecting a command c_i as long the random variable x_i fell within its corresponding probability bin, where $x_i = (i - 1)/m + x$ for $i \in [1, m]$, and $x \sim U(0, 1/m)$.

Mapping Commands to Goals. Given a sampled command $c = ([a, \dots], \{p, \dots\}, \{(v, \tau), \dots\})$, we mapped each action a to a set of goal predicates $\{p_a, \dots\}$ that it achieves. By combining these predicates with the qualifier predicates $\{p, \dots\}$, then adding an existential quantifier for any variables, we created a PDDL goal formula g_c that could be used by a planning algorithm. For example, the command (handover robot human ?key1) where (iscolor ?key1 blue) would be converted to the goal formula (exists (?key1 - key) (and (pickedup-by robot ?key1) (has human ?key1) (iscolor ?key1 blue))).

Some commands had multiple possible groundings. In mDKG (but not in VH), these groundings corresponded to plans with different lengths. As such, we implemented both a *naive* and *efficient* form of literal assistance. In naive assistance, we generated plans for all possible groundings of the lifted command c , and averaged performance metrics across these plans. For example, if key2 was blue, then a grounding of (handover robot human ?key1) where (iscolor ?key1 blue) might be (handover robot human key2), which translates to the the goal (and (pickedup-by robot key2) (has human key2)). In efficient assistance, we simply generated a plan for existential goal formula implied by the lifted command, automatically finding the most efficient way to satisfy that formula.

Plan Generation. We generated an assistive plan in two phases. In the first phase, we generated an optimal *joint* plan to the commanded goal g_c using A^* search with a node budget of 2^{16} (for commands with handover actions, this might involve both agents walking towards a convenient meeting point). If finding a joint plan was unsuccessful, we fell back to computing a plan where only the assistant was allowed to move. In the second phase, we then enforced that only the principal was allowed to move (in mDKG) or pick up items (in VH), and computed the remaining plan to achieve the principal’s true goal g (again, using A^* search with a node budget of 2^{16}).

B.3 Multimodal LLM (GPT-4V)

We used the public preview of GPT-4V (`gpt-4-1106-vision-preview`) as a baseline model in our mDKG experiments. We used zero-shot prompting, providing the model with effectively the same instructions as given to the human raters, along with the final frame of the GIF animation shown to humans (Figure B1). Responses were generated with a temperature of 0.

You’re watching a human and a robot collaborating on a treasure game shown below. The human player wants to collect exactly one of the four target gems. The two players cooperate to achieve the human’s goal.

During the game, the human may communicate with the robot to coordinate their actions. Any communication that happens is shown below the visuals.

The rules of the game are as follows:

- The human player’s goal is to collect one target gem.
- The robot plays an assistive role, and cannot pick up gems on its own.
- The players can move on the white squares.
- The players have a full view of the map at all times.
- Keys are used to unlock doors with the same color (e.g. red keys only unlock red doors).
- Each key can only be used once. Keys disappear after being used.
- One player can pass keys to the other if they’re next to each other.
- Players can occupy the same square and walk past each other.
- The keys and doors are labeled. The labels are shown on the top right corner of each cell.

stimulus.png

In the scenario above, the human player took several actions, then gave the instruction: "Can you hand me that key?"

Given this information, what gem is most likely to be the human agent’s goal? Which key(s) should the robot pick up to best assist the human?

Please respond in the following JSON format:

```
{ "goal": [red/yellow/blue/green], "assist": [key1, key2, etc.] }
```

Figure B1: GPT-4V zero-shot prompt for an example key assistance problems in Doors, Keys & Gems.

C HUMAN EXPERIMENTS

Here we provide additional details about our human experiments for the mDKG domain. Participants first completed a tutorial that explained the task and experimental interface, then answered 5 comprehension questions before proceeding to the main experiment. In the main experiment, they were shown 15 stimuli in randomized order, which involved either keys or doors as assistance options (Figure C1).

To incentivize accurate but calibrated responses, participants were rewarded for accurately guessing the true goal or selecting the optimal assistance option. Specifically, they earned $1/N$ bonus points if they selected N goals out of which one was the true goal, but 0 points if none of their selected goals was the true goal. They earned one additional point if they selected exactly the right assistance options. Participants were paid US\$1 for every 30 bonus points they earned, on top of a base pay of US\$16/hr (US\$4 per experiment).

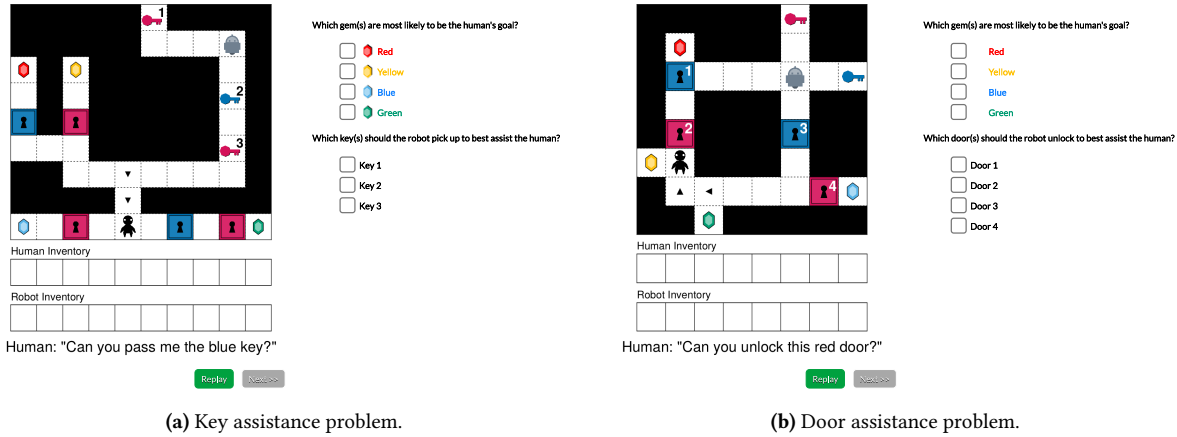


Figure C1: Web interface for human experiments. Participants had to select checkboxes for all gems they thought were likely to be human principal’s goal, and select all keys or doors they thought the robot assistant should pick-up or unlock.

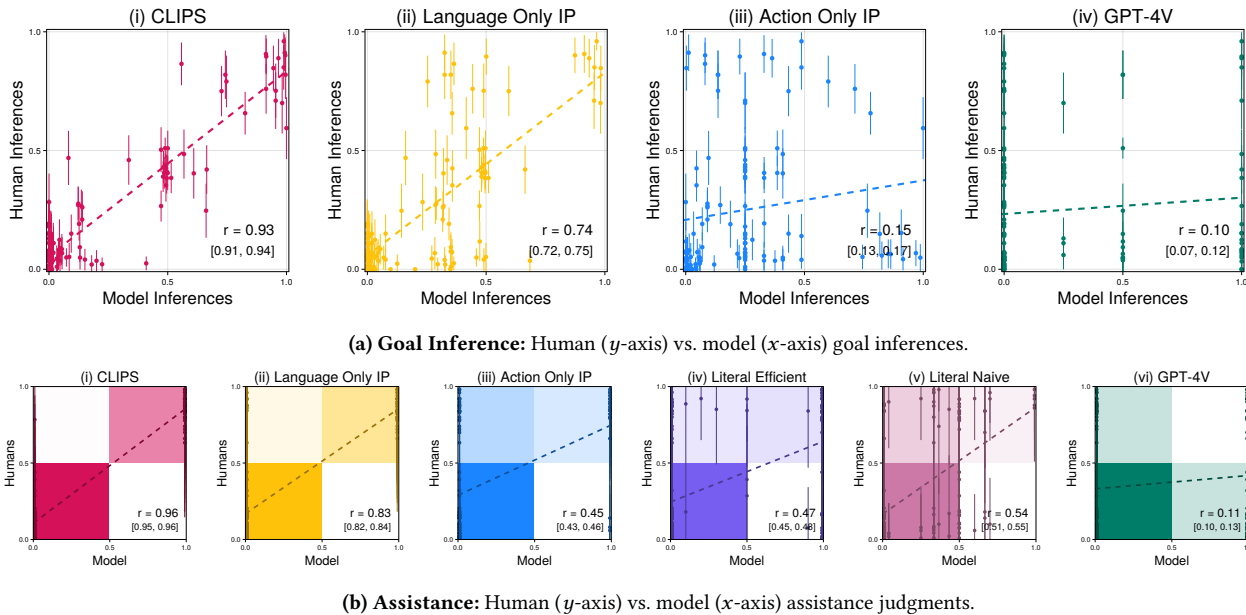


Figure C2: Human-model correlation plots for Doors, Keys & Gems. In (a), we show scatter plots comparing average human goal inferences vs. model goal inferences across all stimuli. In (b), we compare average human assistance judgments vs. model outputs by overlaying a scatter plot over a binarized confusion matrix. Error bars correspond to 95% confidence intervals for the mean human response. At the bottom-right of each plot we report Pearson’s r (95% bootstrap CI below).

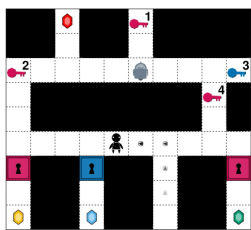
In Figure C2, we show correlation plots that provide a more comprehensive view of how well each model matches human responses. As Figure C2a shows, CLIPS best reproduces human goal inferences. In contrast, unimodal inverse planning tends to be under-confident relative to humans, since they only incorporate evidence from one modality rather than two. GPT-4V tends to be over-confident on stimuli where humans are not (and vice versa), resulting in a very low correlation overall.

To compare human assistance judgments with model outputs, Figure C2b shows scatter plots laid over binarized confusion matrices, generated by thresholding human and model responses at 0.5, then plotting bin counts as a heatmap. Since many of our models return either 1 or 0 for assistance, this provides a better way to visualize human-model agreement. CLIPS best matches human judgments, followed by language only inverse planning. The other models often fail to select an assistance option in cases where humans do (deeper colored top-left quadrants). GPT-4V selects options in many cases where humans do not (bottom-left quadrant), achieving the lowest correlation.

D GPT-4V RESPONSES

To better illustrate the capabilities and limitations of multimodal LLMs, we show several example responses by GPT-4V in Figure D1.

(a) Ambiguous predicates.



Human Inventory
 Robot Inventory

Human: "Can you pass me the red key?"

GPT4V Response:

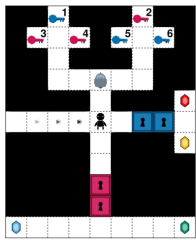
Based on the image and the rules provided, the human player is asking for the red key, which suggests that their goal is to collect the red gem. The red gem is behind the red door labeled "1," and the red key is required to unlock it. Since the human player can only collect one target gem, and they are asking for the red key, it is most likely that the red gem is their goal.

The robot should pick up the red key to assist the human player. The robot cannot pick up gems, so it only needs to focus on the keys that will help the human player achieve their goal.

Here is the response in the requested JSON format:

```
{ "goal": ["red"], "assist": ["key1"] }
```

(b) Ambiguous number.



Human Inventory
 Robot Inventory

Human: "Can you pass me the red keys?"

GPT4V Response:

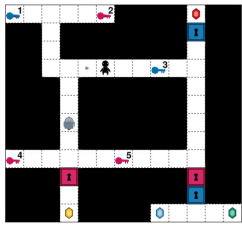
Based on the image and the information provided, the human player has asked the robot to pass the red keys. This implies that the human player is aiming for a gem that requires passing through a red door. In the image, the red gem is behind a red door, so it is the most likely target for the human player.

The robot should pick up the keys that will allow the human to reach the red gem. The red gem is behind door 3, which requires key 3. However, to reach key 3, the robot first needs to unlock door 1 with key 1. Therefore, the robot should pick up key 1 and then key 3 to assist the human player.

Here is the response in the requested JSON format:

```
{ "goal": ["red"], "assist": ["key1", "key3"] }
```

(c) Joint instructions.



Human Inventory
 Robot Inventory

Human: "I'll get the blue key. Can you pick up a red key?"

GPT4V Response:

Based on the information provided and the layout of the game board, the human player has expressed an intention to pick up the blue key. This key can unlock the blue door, which is blocking the path to the green gem. Since the human player's goal is to collect one target gem and they are going for the blue key, it is reasonable to assume that the green gem is their intended target.

The robot, in an assistive role, has been asked to pick up a red key. There are two red keys (key4 and key5) and two red doors on the map. However, since the human is going for the green gem, the robot should pick up the red key that will assist in this goal. Key4 is the red key that unlocks the red door on the path to the green gem. Therefore, the robot should pick up key4.

Here is the JSON response with the likely goal and the key the robot should pick up:

```
{ "goal": ["green"], "assist": ["key4"] }
```

Figure D1: Example responses by GPT-4V. GPT-4V produces two paragraphs of reasoning text, then gives its final response in JSON.

As can be seen from Figure D1, GPT-4V is highly inaccurate in its ability to perceive objects in the provided image and determine their spatial relationships, sometimes confabulating object labels that do not exist (Fig. D1a–b), failing to notice which gems are locked behind which doors (Fig. D1b–c), and failing to perceive which items are along a path to the goal (Fig. D1c). This weakness at spatial perception makes current multimodal LLMs inadequate for assistive tasks like ours, which require coherent spatial reasoning.

However, the reasoning text produced by GPT-4V appears to demonstrate internal means-ends consistency (keys are required to unlock doors of the same color, which in turn allow the human to reach the gems behind them). As prior work has shown [60], this suggests that if an abstract symbolic representation is provided to a state-of-the-art LLM with the appropriate high-level relations (spatial connectivity etc.), they may be able to perform significantly better. Nonetheless, deriving an abstract symbolic representation from the underlying gridworld is itself a form of complex spatial reasoning [5], which likely presents a challenge to contemporary LLMs.

A full set of GPT-4V’s responses can be found in the Supplementary Information (<https://osf.io/v8ru7/>).

E ALTERNATE ASSISTANCE METHODS

In the main text, we reported results for an *offline* version of the Q_{MDP} assistance policy in Algorithm 2, which keeps the goal posterior fixed after observing $a_{1:T}^h$ and $u_{1:T}$, rather than continuing to update the posterior as the human principal takes more actions. Here we report results for the *online* version of Q_{MDP} assistance. We also report results for an alternative assistance method, $\bar{\pi}_{MDP}$, which is constructed by averaging over policies π for each inferred goal g : Rather than minimizing the expected Q -value as in Q_{MDP} assistance [45], $\bar{\pi}_{MDP}$ first samples a goal g from the posterior, then follows the policy π for that goal. This can be viewed as Thompson sampling [72] over policies.

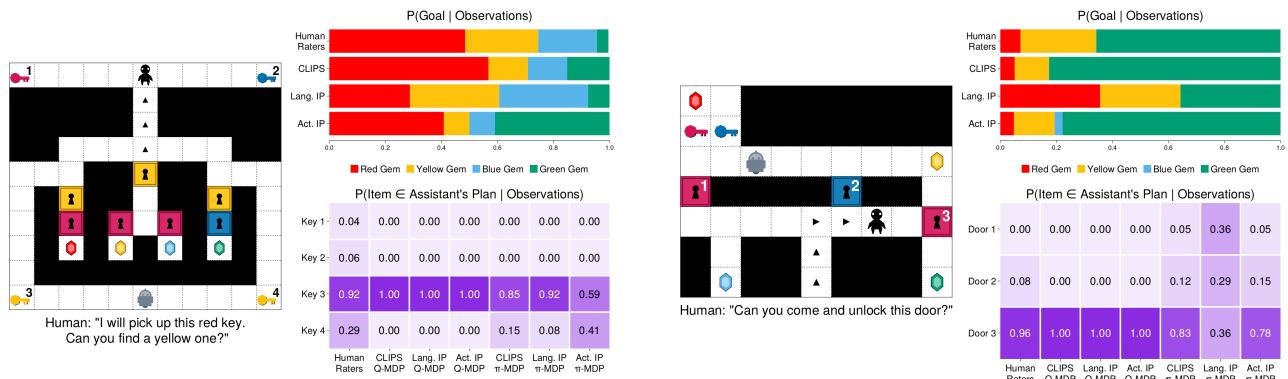


Figure E1: Examples comparing (offline) Q_{MDP} and $\bar{\pi}_{MDP}$ assistance in Doors, Keys & Gems. In Q_{MDP} assistance, we repeatedly take the best action that minimizes the expected goal achievement cost with respect to the posterior over goal specifications. In $\bar{\pi}_{MDP}$ assistance, we randomly select a single goal specification g from the posterior, and follow the optimal joint policy for that goal.

Table E1: Comparison of assistance methods measured in terms accuracy (precision and recall for assistance options), helpfulness (plan length and human cost relative to CLIPS with offline Q_{MDP} assistance), and human similarity (correlation of assistance options with mean human ratings) Metrics are averaged across the dataset per domain, with standard errors reported in brackets.

Inference	Assistance	Accuracy		Helpfulness		Human Sim.
		Assist. Prec.	Assist. Rec.	Rel. Plan Length	Rel. Human Cost	Assist. Cor.
<i>Doors, Keys & Gems</i>						
CLIPS	Q_{MDP} (Offline)	0.97 (0.03)	0.97 (0.03)	1.00 (0.00)	1.00 (0.00)	0.96 (0.01)
Lang. Only IP	Q_{MDP} (Offline)	0.90 (0.06)	0.83 (0.06)	1.26 (0.10)	1.18 (0.07)	0.83 (0.01)
Action Only IP	Q_{MDP} (Offline)	0.43 (0.09)	0.40 (0.09)	1.68 (0.15)	1.46 (0.08)	0.45 (0.01)
CLIPS	Q_{MDP} (Online)	1.00 (0.00)	1.00 (0.00)	0.99 (0.01)	0.99 (0.00)	0.97 (0.00)
Lang. Only IP	Q_{MDP} (Online)	0.90 (0.06)	0.83 (0.06)	1.26 (0.10)	1.18 (0.07)	0.83 (0.01)
Action Only IP	Q_{MDP} (Online)	0.88 (0.06)	0.90 (0.06)	1.15 (0.07)	1.09 (0.03)	0.83 (0.01)
CLIPS	$\bar{\pi}_{MDP}$	0.91 (0.03)	0.91 (0.03)	1.09 (0.04)	1.08 (0.02)	0.96 (0.00)
Lang. Only IP	$\bar{\pi}_{MDP}$	0.79 (0.05)	0.77 (0.05)	1.31 (0.09)	1.19 (0.04)	0.85 (0.01)
Action Only IP	$\bar{\pi}_{MDP}$	0.72 (0.04)	0.41 (0.05)	1.58 (0.12)	1.45 (0.05)	0.58 (0.01)
<i>VirtualHome</i>						
CLIPS	Q_{MDP} (Offline)	0.87 (0.04)	1.00 (0.00)	1.00 (0.00)	1.00 (0.00)	–
Lang. Only IP	Q_{MDP} (Offline)	0.59 (0.05)	0.96 (0.03)	1.33 (0.07)	1.35 (0.08)	–
Action Only IP	Q_{MDP} (Offline)	0.61 (0.07)	0.84 (0.07)	1.30 (0.07)	1.33 (0.07)	–
CLIPS	Q_{MDP} (Online)	0.90 (0.04)	1.00 (0.00)	1.12 (0.15)	1.13 (0.17)	–
Lang. Only IP	Q_{MDP} (Online)	0.59 (0.05)	0.96 (0.03)	1.33 (0.07)	1.35 (0.08)	–
Action Only IP	Q_{MDP} (Online)	0.84 (0.06)	0.92 (0.06)	1.03 (0.04)	1.04 (0.05)	–
CLIPS	$\bar{\pi}_{MDP}$	0.59 (0.06)	0.65 (0.06)	1.22 (0.06)	1.33 (0.07)	–
Lang. Only IP	$\bar{\pi}_{MDP}$	0.45 (0.05)	0.57 (0.06)	1.37 (0.08)	1.48 (0.09)	–
Action Only IP	$\bar{\pi}_{MDP}$	0.46 (0.06)	0.45 (0.05)	1.61 (0.14)	1.77 (0.16)	–

Qualitative analysis. Figure E1 illustrates how (offline) Q_{MDP} assistance differs from $\bar{\pi}_{MDP}$. While $\bar{\pi}_{MDP}$ assistance better reproduces the gradedness seen in human assistance judgments, this comes at the cost of reduced safety and helpfulness. By following a policy π to goal g according to how likely g is, $\bar{\pi}_{MDP}$ fails to aggregate information across multiple uncertain goals when deciding the best action to take. This is most clearly illustrated in Figure E1c, where CLIPS assigns higher probability to the blue gem over the green gem. Q_{MDP} assistance safely unlocks *Door 3* instead of *Door 2*, since doing so will ensure that both the blue gem and green gem are reachable. In contrast, $\bar{\pi}_{MDP}$ unlocks *Door 2* on average 66% of the time, which is unsafe. However, as Figure E1d shows, $\bar{\pi}_{MDP}$ is slightly better when CLIPS produces the wrong goal inferences (due to the sub-optimality of the human’s actions with respect to the true goal). In this single failure case, Q_{MDP} confidently executes the wrong actions, whereas $\bar{\pi}_{MDP}$ assistance is sometimes successful.

Quantitative analysis. We present the performance of different assistance methods in Table E1. Compared to offline Q_{MDP} assistance, $\bar{\pi}_{MDP}$ assistance was less accurate, and produced longer plans on average. This highlights the importance of minimizing expected cost under uncertainty, rather than just sampling a likely goal to follow. As might be expected, online Q_{MDP} assistance was more accurate, since the posterior over goals was updated as new actions were observed. This significantly improved the performance of action-only inverse planning (note that language-only inverse planning is the same regardless of assistance method, since the human principal communicates no further instructions). However, multimodal goal inference via CLIPS combined with online Q_{MDP} assistance was still the most accurate, since it leveraged linguistic information from the beginning of assistance.

Apart from one failure case in VirtualHome where the assistant and simulated human principal were caught in a livelock, CLIPS with online Q_{MDP} assistance was slightly more helpful on average than offline assistance. Nonetheless, offline assistance was almost as helpful. This suggests that once enough actions are observed to disambiguate the principal’s instruction, CLIPS can still assist effectively even if it has partial to no observability of the principal’s subsequent actions.