

The Neuro-Symbolic Inverse Planning Engine

Author Names Omitted for Anonymous Review. Paper-ID [add your ID here]

Abstract—Human beings are social creatures. We routinely reason about other agents, and a crucial component of this social reasoning is inferring people’s goals as we learn about their actions. In many settings, we can perform intuitive but reliable goal inference from language descriptions of agents, actions, and the background environments. In this paper, we study this process of language driving and influencing social reasoning in a probabilistic goal inference domain. We propose a neuro-symbolic model that carries out goal inference from linguistic inputs of agent scenarios. The “neuro” part is a large language model (LLM) that translates language descriptions to code representations, and the “symbolic” part is a Bayesian inverse planning engine. To test our model, we design and run a human experiment on a linguistic goal inference task. Our model closely matches human response patterns and better predicts human judgements than using an LLM alone.

I. INTRODUCTION

Humans are deeply social creatures, finely attuned to the inner lives of other people (Wellman [29], Saxe [21], Tomasello [25]). We not only think, but think routinely *about each other*. This capacity for *social reasoning*, broadly, is *abstract* and relies on a fine grained knowledge of distinct but highly *interrelated* mental states within others.

This paper considers how **language can inform social reasoning**, with a particular focus on **understanding language about people’s actions, plans, and goals**. Language is an especially powerful means of conveying the abstract, interrelated nature of concepts about agents in the world—we can talk about someone’s actions (*I saw Annie heading into the living room*), directly convey their mental states (*Annie desperately wants a toy in there*), describe relevant but abstract aspects of the world in which they act (*that cupboard Annie is trying to open is locked inside*), or even pose the questions we want others to reason about (*which toy do you think she wants the most?*). Language about any one of these aspects informs how we might reason about someone else in holistic ways, changing our understanding and downstream predictions about what they might know, want, or do.

This paper focuses then on **understanding language that describes inverse planning and goal inference problems**—how language can convey information about the world and the actions within it, in order to drive inferences about unseen goals motivating these actions. We make the following contributions:

- 1) We introduce a new model, the *Neuro-symbolic Inverse Planning Engine* (NIPE) that leverages LLMs to *translate language into program expressions* that compositionally construct and condition a *probabilistic generative model over agents and environments*, then

leverages *Bayesian inverse planning* to perform goal inference over this symbolic representation.

- 2) To evaluate this model, we compare our model against *human social reasoning on a linguistic inverse planning domain* inspired by prior cognitive work in social cognition Zhi-Xuan et al. [31]. We design this domain to require jointly reasoning about agent actions and varying information about the world itself.
- 3) We find that our model *well-predicts human judgments* across a variety of distinct social reasoning scenarios. We also find that our model *outperforms LLM baselines used to directly perform the same goal inference tasks*. We find that the largest LLM alone (GPT-4 [19]) correlates moderately well with certain human judgments on this domain, particularly on simpler scenarios. We discuss trends between LLMs of different sizes and on different task variants in our stimuli.

II. LINGUISTIC INVERSE PLANNING TASKS AND HUMAN EXPERIMENT

We begin by describing the domain of tasks and human experiment we use throughout this paper, to provide intuition for the modeling experiments that follow. We construct a domain of **linguistic inverse planning stimuli** designed to evaluate how humans and models **infer the goals of agents acting within varied environments**. We design stimuli to convey varying information about the *world configuration and dynamics* in which agents behave, as well as the *observed actions of agents* with unknown goals.

Our tasks are inspired by the spatial navigation stimuli used in prior social reasoning experiments (Zhi-Xuan et al. [31], Alanqary et al. [1]). Our experiments adapt this basic domain into a set of linguistic inverse planning stimuli, and introduce a range of variations designed to test distinct world dynamics, spatial information, and vague or incomplete information about agent actions conveyed linguistically.

A. Planning and inverse planning concepts in the Gameshow Domain

Each stimuli in our experiment (table I) begins with a description of a “*gameshow*” scenario in which a player is trying to reach one of several goal *trophies* placed around an obstacle course. This obstacle course also includes *locked doors* that can only be unlocked with keys placed around the course, and which additionally require multi-step reasoning about an agent’s longer term plans. Based on each scenario description, which includes varying information about the particular obstacle course and the actions of an agent within it, we pose a linguistic goal inference query: *which of the*

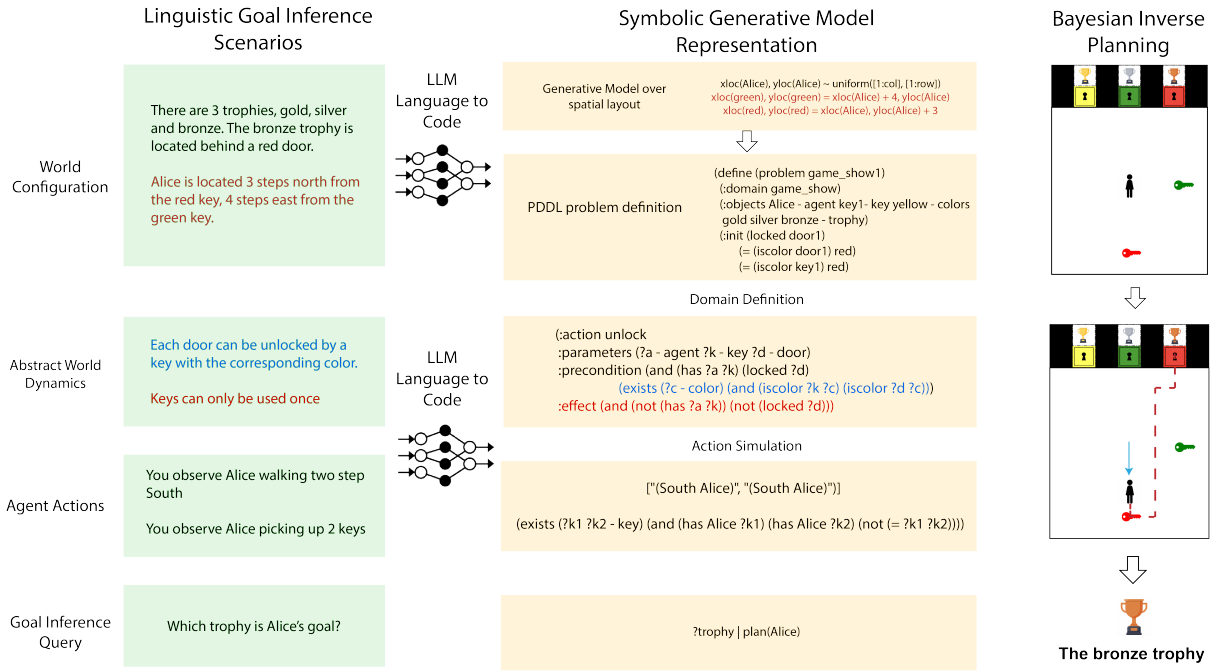


Fig. 1. Overview of our model, the *Neuro-Symbolic Inverse Planning Engine* (NIPE), which reasons about linguistic goal inference scenarios. (Left, green) We design and measure human judgments on a *linguistic goal inference* domain. Our model works by first using a LLM to *translate* language into symbolic program semantics that collectively condition a probabilistic generative model over possible environments and agents (Middle, yellow), and answers queries via *Bayesian inverse planning* (Right) in the conditioned probabilistic model, inferring likely goals using a nested planner to model the agent.

Linguistic planning and inverse planning concepts across tasks	
World configuration	<i>There are three trophies placed through this obstacle course: gold, silver, and bronze; The gold trophy is currently 5 steps East of Alice; The silver trophy is behind 1 door.</i>
World dynamics	<i>On this course, each lock can only be opened by a key that is the same color as that lock; In this game, players lose points if they pick up more keys than needed.</i>
Agent actions	<i>You see Alice picking up a green key; You see Alice walking a step towards the South In this game, players lose points if they pick up more keys than needed.</i>
Goal inference	<i>Considering each of the possible trophies on this course, how likely is it that Alice is trying to reach this trophy?</i>

TABLE I
EXAMPLES FROM THE LINGUISTIC “GAMESHOW” GOAL INFERENCE DOMAIN. EACH SCENARIO’S LANGUAGE SPANS A VARIETY OF CONCEPTS RELEVANT TO THE FINAL GOAL-INFERENCE QUERY, FROM INFORMATION ABOUT THE ENVIRONMENT CONFIGURATION TO AN AGENT’S ACTIONS.

possible trophies do you think the player trying to reach? Full details on the stimuli are included in the Appendix.

B. Human Experiment

We design a human experiment and collect participants’ judgments on these goal inference tasks (see Appendix).

III. OUR MODEL: THE NEURO-SYMBOLIC INVERSE PLANNING ENGINE (NIPE)

A. Language-to-code translation using LLMs

In this section, we describe a general framework for conditioning on the diverse range of linguistic information in these stimuli—varying information about the spatial environment, world dynamics, and agent actions—that supports inferences

about agents based on a principled, probabilistic framework for Bayesian inverse planning.

In our study, this is implemented by using LLMs to *translate sentences into code*, which compositionally construct a symbolic model over an extended, multi-sentence linguistic scenario (see fig. 1, *LLM Language to Code*).

Our translation approach leverages the few-shot prompting capabilities of LLMs trained jointly on language and code. Models were sampled with *temperature*=1.2. To ensure translation robustness, we rejection sample translations to ensure that translated code was syntactically valid and executable.

First, we prompt LLM to translate sentences about the environment (World Configuration) into a code `CONDITION` that describes this arrangement. We then sample concrete grid-

world maps from a generative model over maps that satisfy this. To ensure that maps are valid for our environment, we define and condition a relatively restricted generative model over possible maps and trophy placements that ensures that these basic validity conditions are met.

Each sampled map from this generative model is a concrete PDDL scene representation that describes the locations of the initial agent, objects, and their attributes.

Sentences about domain rules (Abstract World Dynamics) are straightforwardly translated into operator definitions expressed in PDDL, which define an action and its necessary parameters, preconditions for trying this action, and effects on the environment. Collectively, the scene representation and the domain rules define a complete, concrete planning domain defined in PDDL.

Finally, sentences that describe agents’ actions, such as “You observe Alice walking 3 steps East.”, can be directly translated to a low-level action sequence with LLM.

B. Goal Inference through Bayesian Inverse Planning

Once we have an initial PDDL environment representation s_0 and a sequence of agent actions $a_{1:t}$, we can perform Bayesian inverse planning to infer a posterior distribution over goals g given actions $a_{1:t}$ and the initial state s , which is proportional to the joint distribution over goals and actions:

$$P(g|s_0, a_{1:t}) \propto P(g, a_{1:t}|s_0) = P(g|s_0)P(a_{1:t}|g)$$

Here $P(g|s_0)$ defines the prior distribution over possible goals in the initial state s_0 , and $P(a_{1:t}|g)$ is a model of how the agent selects actions $a_{1:t}$ given a particular goal g . For the goal prior, we model agents as being more likely to pick up goals that are closer to them. To model actions, we assume $P(a_{1:t}|g)$ decomposes into the probability of selecting an action a_t at each state s_t , given by a Boltzmann-rational policy π :

$$P(a_{1:t}|g) = \prod_{\tau=1}^t \pi(a_\tau|s_\tau, g)$$

$$\pi(a_t|s_t, g) = \frac{\exp \beta Q_g(s_t, a_t)}{\sum_{a'_t} \exp \beta Q_g(s_t, a'_t)}$$

where $Q_g(s_t, a_t)$ denotes the (negative) cost of the reaching the goal g by taking action a_t from s_t , and β is a rationality parameter, also called the inverse temperature. To compute Q_g for each state s_t and action a_t encountered during inference, we use A* search as an optimal-cost planner in the PDDL domain to determine the cost of the optimal plan.

With these components of the model specified, goal inference can be performed exactly by computing the joint probability $P(g, a_{1:t}|s_0)$ of each possible goal g with the observed sequence of actions $a_{1:t}$, which can then be normalized to obtain the posterior distribution $P(g|s_0, a_{1:t})$.

IV. MODEL EXPERIMENTS

We compare our model with human judgements on the gameshow domain of social reasoning tasks. We also evaluate two LLM baselines to assess the value of an explicit Bayesian inverse planning model.

For our model, we simulate model answers to each goal inference stimulus (*which trophy will Alice pick up?*), outputting a *posterior probability* over all of the possible trophies. We also conduct two baseline experiments using two different LLM models (GPT-3.5 [20] and GPT-4 (gpt-4-0314) [19]) to directly provide goal inference judgments for each stimulus, with no external inference or planning engine. We assess LLMs to exactly replicate the conditions of the human experiment with zero-shot prompting. To account for answer variability, we sample $n = 30$ answers for each stimulus, rejecting any answers that are not in the correct answer format.

V. RESULTS AND DISCUSSION

We evaluate our model and baselines in comparison to human judgments, both across all of the stimuli and within each task variant category. We find that:

a) *Our model (NIPE) correlates strongly overall with human judgments across the goal inference experiment.*: We calculate correlations between the *posterior probability judgments* produced by our model and the mean per stimulus judgments (Likert scale ratings) across human subjects, and find that our model is strongly correlated with human judgments in the experiment overall (*Our Model*, $R=0.927$, 95% CI=(0.877, 0.953)). All 95% CI are bootstrap sampled confidence intervals each with 1000 samples.

b) *The largest LLM baseline (GPT-4) correlates moderately with human judgments overall; GPT3.5 correlates poorly.*: We additionally calculate correlations between LLMs prompted directly to produce goal inference judgments (mean over sampled LLM judgements for each stimulus, and the mean over human judgments for each stimulus across participants). In comparison to our model, we find that GPT-4 correlates well, but more weakly, with human judgments, *GPT-4 Rating*, $R=0.658$, CI=(0.403, 0.79)) overall, and GPT-3.5 correlates much more poorly with human judgments overall, (*GPT-3.5 Rating*, $R=0.100$, CI=(-0.145, 0.361)). These aggregate trends are better understood by examining correlations between human judgments within each task variant, as we discuss next.

c) *NIPE correlates strongly with human judgements within each task variant; GPT-3.5 and GPT-4 are less robust to complex world dynamics and spatial information.*: Inspecting trends *within* each individual task variant (fig. 2, columns correspond to *task variants*, and rows correspond to *models*) provides insight into when models deviate from human judgements.

NIPE consistently correlates with human judgments across *all* task variants, from the simpler *Generic* stimuli to the more complex *Colored* and *Spatial* stimuli (fig. 2, bottom). In contrast, on the more complex tasks—particularly the *Colored-Different* variant, which specifies a new, counterintuitive rule particular to this domain (*doors can only be opened by keys that are a different color from the door*), and the *Spatial* variant, which requires considering how an efficient, path-planning agent would navigate towards their goal—the LLM

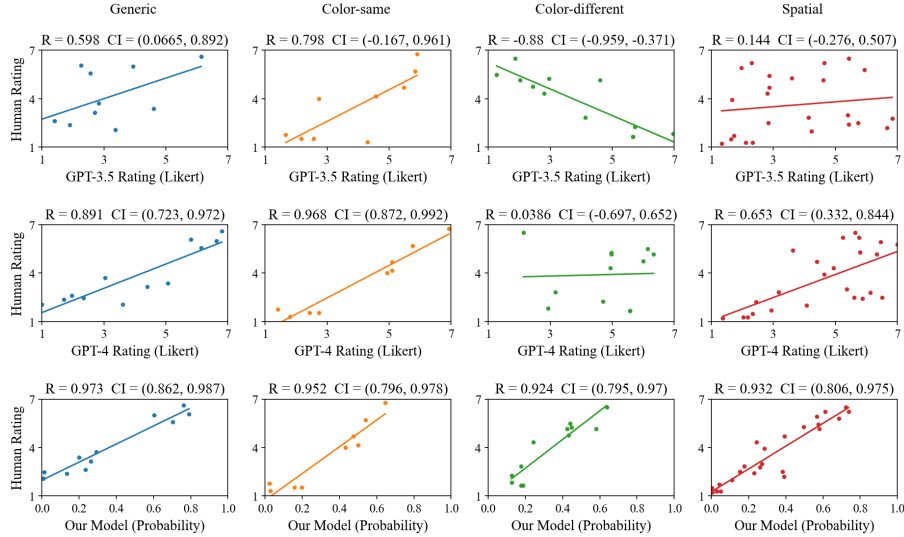


Fig. 2. Correlations within specific task variants highlight more granular trends in model behavior. Our model (*Bottom*) correlates strongly with human judgements across all four task variants. In comparison, the LLM baselines generally show poorer performance on the more complex variants (the *Color-Different* variant, which introduces a new, counterintuitive rule about the world dynamics; and the *Spatial* variant, which requires reasoning specifically about agents as efficient path planners).

baselines show revealingly different performance in comparison to human reasoning.

As seen in fig. 2, GPT-3.5 actually correlates *negatively* with human judgments on the *Colored-Different*—it both ignores the stated rule about the different colored keys and appears to make the opposite, incorrect assumption, perhaps based on simpler color-based pattern matching; and shows nearly zero correlation with human judgments on the explicit *Spatial* path-planning variant. In comparison, GPT-4 correlates well with human judgments on both of the simpler task variants, and shows a moderate (but positive) correlation on the *Spatial* variant (though we discuss stimuli-specific failure modes next, suggesting that the model may still be using simpler distance heuristics, rather than robustly considering an efficient path-planning agent.) On the *Colored-Different* domain, GPT-4 has nearly zero-correlation, suggesting that it is highly uneven in taking the specific rule about key colors into account—as shown in fig. 2, it predicts human judgments well in some cases and in others, again produces nearly the opposite overall judgment, suggesting it has ignored the rule.

d) Humans make fine-grained assumptions about agents as rational planners; LLMs do not always do so.: Our model explicitly reasons about agents as rational, goal-directed actors and planners. Do LLMs do so as well? Qualitative analysis of individual stimuli highlights subtle but important discrepancies between human assumptions and LLM-based goal inferences on the most complex *Spatial* stimuli. We highlight a representative example here:

There are three trophies placed throughout this obstacle course: gold, silver, or bronze. Alice is a participant on this course. The gold trophy is currently 2 steps East of Alice. The silver trophy is 3 steps South from her. The bronze trophy is 5 steps

South from her. You see Alice walking 4 steps South.

Both humans and our model strongly predict that Alice is headed towards the *bronze* trophy on this domain. In contrast, both the LLM baselines predict most strongly that Alice is heading towards the *silver trophy* (perhaps because it is closer to Alice), though in fact, humans tend to agree that because she has walked *past* this trophy, as an efficient actor, this is likely not her goal.

VI. CONCLUSION AND FUTURE DIRECTIONS

We conclude with several important avenues for future work. First, while we evaluate all LLMs using a simple prompting regime, future work should evaluate whether these results are consistent or improve with other prompting and evaluation methods such as chain-of-thought [28, 13], particularly on complex tasks such as those that invoke efficient spatial path planning, which prior work also finds to be an outstanding challenge for LLMs [8].

More generally, we design the gameshow domain to permit future extensions to describe a much broader range of social reasoning tasks, many of which have been studied in the cognitive and AI literature and that we might in turn express linguistically, including: scenarios that specifically invoke agent belief; involve boundedly rational agents; and even involve multiple agents collaborating or competing. Future work can also extend this domain to *multimodal stimuli* that include both linguistic information *and* demonstrations or images, e.g., to convey complex agent trajectories or environment layouts. Our model suggests a general framework for extending our language-to-code approach and Bayesian inverse planning model to address these more complex inferences, building on other computational and cognitive models that reason jointly over agent beliefs, goals, and multiple agents.

REFERENCES

- [1] Arwa Alanqary, Gloria Z Lin, Joie Le, Tan Zhi-Xuan, Vikash K Mansinghka, and Joshua B Tenenbaum. Modeling the mistakes of boundedly rational agents within a Bayesian theory of mind. *arXiv preprint arXiv:2106.13249*, 2021.
- [2] Renée Baillargeon, Rose M Scott, and Lin Bian. Psychological reasoning in infancy. *Annual review of psychology*, 67:159–186, 2016.
- [3] Chris L Baker, Rebecca Saxe, and Joshua B Tenenbaum. Action understanding as inverse planning. *Cognition*, 113(3):329–349, 2009.
- [4] Chris L Baker, Julian Jara-Ettinger, Rebecca Saxe, and Joshua B Tenenbaum. Rational quantitative attribution of beliefs, desires and percepts in human mentalizing. *Nature Human Behaviour*, 1(4):0064, 2017.
- [5] Gagan Bansal, Besmira Nushi, Ece Kamar, Eric Horvitz, and Daniel S Weld. Is the most accurate ai the best teammate? optimizing ai for teamwork. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 11405–11414, 2021.
- [6] Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, Harsha Nori, Hamid Palangi, Marco Tulio Ribeiro, and Yi Zhang. Sparks of Artificial General Intelligence: Early experiments with GPT-4, 2023.
- [7] Micah Carroll, Rohin Shah, Mark K Ho, Tom Griffiths, Sanjit Seshia, Pieter Abbeel, and Anca Dragan. On the utility of learning about humans for human-ai coordination. *Advances in neural information processing systems*, 32, 2019.
- [8] Anthony G Cohn and Jose Hernandez-Orallo. Dialectical language model evaluation: An initial appraisal of the commonsense spatial reasoning abilities of LLMs. *arXiv e-prints*, pages arXiv–2304, 2023.
- [9] Katherine M Collins, Catherine Wong, Jiahai Feng, Megan Wei, and Josh Tenenbaum. Structured, flexible, and robust: benchmarking and improving large language models towards more human-like behavior in out-of-distribution reasoning tasks. In *Proceedings of the Annual Meeting of the Cognitive Science Society*, volume 44, 2022.
- [10] Gergely Csibra, György Gergely, Szilvia Biró, Orsolya Koos, and Margaret Brockbank. Goal attribution without agency cues: the perception of ‘pure reason’ in infancy. *Cognition*, 72(3):237–267, 1999.
- [11] Daniel Clement Dennett. *The Intentional Stance*. MIT Press, 1981.
- [12] Julian Jara-Ettinger, Hyowon Gweon, Laura E Schulz, and Joshua B Tenenbaum. The naïve utility calculus: Computational principles underlying commonsense psychology. *Trends in cognitive sciences*, 20(8):589–604, 2016.
- [13] Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large language models are zero-shot reasoners. In *Advances in Neural Information Processing Systems*, 2022.
- [14] Michal Kosinski. Theory of mind may have spontaneously emerged in large language models. *arXiv preprint arXiv:2302.02083*, 2023.
- [15] Bosen Lian, Wenqian Xue, Frank L Lewis, and Tianyou Chai. Online inverse reinforcement learning for nonlinear systems with adversarial attacks. *International Journal of Robust and Nonlinear Control*, 31(14):6646–6667, 2021.
- [16] Bo Liu, Yuqian Jiang, Xiaohan Zhang, Qiang Liu, Shiqi Zhang, Joydeep Biswas, and Peter Stone. Llm+ p: Empowering large language models with optimal planning proficiency. *arXiv preprint arXiv:2304.11477*, 2023.
- [17] Qing Lyu, Shreya Havaldar, Adam Stein, Li Zhang, Delip Rao, Eric Wong, Marianna Apidianaki, and Chris Callison-Burch. Faithful chain-of-thought reasoning. *arXiv preprint arXiv:2301.13379*, 2023.
- [18] Hussein Mozannar and David Sontag. Consistent Estimators for Learning to Defer to an Expert. In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 7076–7087. PMLR, 13–18 Jul 2020. URL <https://proceedings.mlr.press/v119/mozannar20b.html>.
- [19] OpenAI. GPT-4 Technical Report, 2023.
- [20] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744, 2022.
- [21] Rebecca Saxe. Uniquely human social cognition. *Current opinion in neurobiology*, 16(2):235–239, 2006.
- [22] Natalie Shapira, Mosh Levy, Seyed Hossein Alavi, Xuhui Zhou, Yejin Choi, Yoav Goldberg, Maarten Sap, and Vered Shwartz. Clever Hans or Neural Theory of Mind? Stress Testing Social Reasoning in Large Language Models. *arXiv preprint arXiv:2305.14763*, 2023.
- [23] Ben Shneiderman. *Human-centered AI*. Oxford University Press, 2022.
- [24] Tianmin Shu, Abhishek Bhandwaldar, Chuang Gan, Kevin Smith, Shari Liu, Dan Gutfreund, Elizabeth Spelke, Joshua Tenenbaum, and Tomer Ullman. Agent: A benchmark for core psychological reasoning. In *International Conference on Machine Learning*, pages 9614–9625. PMLR, 2021.
- [25] Michael Tomasello. *Origins of human communication*. MIT press, 2010.
- [26] Tomer Ullman. Large Language Models Fail on Trivial Alterations to Theory-of-Mind Tasks. *arXiv preprint arXiv:2302.08399*, 2023.
- [27] Karthik Valmeekam, Alberto Olmo, Sarath Sreedharan, and Subbarao Kambhampati. Large Language Models Still Can’t Plan (A Benchmark for LLMs on Planning and Reasoning about Change). *arXiv preprint*

arXiv:2206.10498, 2022.

- [28] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed H Chi, Quoc V Le, Denny Zhou, et al. Chain-of-Thought Prompting Elicits Reasoning in Large Language Models. In *Advances in Neural Information Processing Systems*, 2022.
- [29] Henry M Wellman. *The child's theory of mind*. The MIT Press, 1992.
- [30] Bryan Wilder, Eric Horvitz, and Ece Kamar. Learning to complement humans. In *Proceedings of the Twenty-Ninth International Conference on International Joint Conferences on Artificial Intelligence*, pages 1526–1533, 2021.
- [31] Tan Zhi-Xuan, Jordyn Mann, Tom Silver, Josh Tenenbaum, and Vikash Mansinghka. Online bayesian goal inference for boundedly rational planning agents. *Advances in neural information processing systems*, 33: 19238–19250, 2020.

A. Additional background

We provide additional context within which to situate our work into how language informs aspects of social reasoning. Empirical evidence suggests that we think intuitively about each other as largely *rational, goal-directed intelligent agents* (Dennett [11], Csibra et al. [10], Baker et al. [3], Baillargeon et al. [2]). Within computational cognitive science and AI, a considerable body of prior work formalizes these intuitions by modeling agents as *rational actors and planners* who choose actions to achieve specified goals or optimize rewards, and in turn models many inferences about mental states as *inverse planning* from observations of actions back to their latent goals, utility functions, or beliefs about world states (Jara-Ettinger et al. [12], Baker et al. [4], Zhi-Xuan et al. [31]). Despite considerable successes in predicting distinct human inferences and modeling cooperative behavior in human-computer interactions (Carroll et al. [7], Mozannar and Sontag [18], Alanqary et al. [1], Bansal et al. [5], Wilder et al. [30], Shneiderman [23]), these approaches—largely *non-linguistic* models that draw inferences from images or demonstrations of agent actions—face scalability and specification challenges inherent to the abstract, complex, and interrelated nature of social reasoning tasks. Modeling agents as goal-directed planners requires implicitly or explicitly representing the *world model* they are planning over, which often must be hand-engineered (as in Baker et al. [4] and Zhi-Xuan et al. [31]) or learned from expensive exploration (as in Lian et al. [15]). Inferring mental states, or even posing the task itself, is challenging when framed as an inductive problem from images and demonstrations alone (as in domains like Shu et al. [24]).

In part due to these challenges, the recent advances made by *large language models* (LLMs) have driven an explosion of interest in how these models can be used to plan (Valmeekam et al. [27]) or draw social inferences (Bubeck et al. [6], Shapira et al. [22]) from language—they suggest a new frontier of cognitive models and AI systems that can leverage language to specify abstract world knowledge, domain-specific world dynamics, and direct information about agents and their mental states. However, early efforts to apply LLMs directly to social reasoning problems often take a highly simplified approach to evaluating their capacity for social reasoning or their “*theory of mind*”, drawing broad conclusions based on a narrow set of templated cognitive tasks (like the famous “Sally Anne tasks”) (Kosinski [14]). Other evaluations caution against drawing general conclusions from such restricted task settings, pointing out that modifications and extensions beyond these simple problem settings can produce counterintuitive results that contradict the holism and generality of these inferences (Shapira et al. [22], Ullman [26]); and several recent evaluations suggest that robust planning in general remains an outstanding challenge for large language models (Valmeekam et al. [27], Bubeck et al. [6]). Several recent approaches seek to address these challenges by augmenting LLMs with interfaces to external modules for representing and simulating complex aspects of world state, particularly using classical AI representations and engines for forward planning (Collins et al. [9], Liu et al. [16], Lyu et al. [17]). Emulating human-like social reasoning and knowledge, however, requires models that can robustly draw *unified inferences about actions, mental states, and the world*, relating plans and beliefs in holistic, flexible, and ideally interpretable and safely deployable way to language; and linguistic evaluations that move in turn towards the richness and breadth of human-like social reasoning, flexible and extensible enough to incorporate these interrelated aspects of human mental states and world knowledge.

B. Full details on stimuli construction

We design each stimulus to span several core conceptual categories relevant to reasoning about agents and their actions in the world (table I, *top*), including:

- **World configuration:** information about goal trophies, keys, and their attributes, as well as their relative locations, at differing levels of specificity.
- **Abstract world dynamics:** causal dynamics within the environment, such as rules determining which keys can open which doors.
- **Agent actions:** descriptions of an agent’s behavior before they reach the goal, like *picking up keys* or *walking* in different directions towards keys or trophies.

C. Task variants

We design four distinct variants of gameshow stimuli for a total of $n=18$ stimuli, which vary the complexity of the specified world configuration, the world dynamics, and the type of observed agent actions:

- **Generic keys:** Goal trophies are placed behind locked doors, and must be unlocked by keys placed around the obstacle course. Any key on the course can be used to open any other door.
- **Same-colored keys:** This variant adds additional complexity in the world dynamics. Keys and doors have colors. Doors can only be unlocked by keys that are the same color as the door.
- **Different-colored keys:** Doors can only be unlocked by keys that are a *different* color from the door.
- **Spatial:** This contains specific spatial descriptions of the obstacle course and agent actions. There are no keys or doors; agents can walk freely towards various goal trophies.

D. Human linguistic inverse planning experiment

We include additional details on our human experiment.

We recruit a total of $n=46$ participants from Prolific¹; each viewed a random batch of 7 stimuli from a specific task variant, and for each stimulus, produced goal inference judgements about each possible goal trophy on a 1-7 Likert scale (*1: definitely not this trophy; 7: definitely this trophy*). Scenarios were presented for a fixed amount of time before an answer was permitted (90 seconds in the *Spatial* condition; 45 seconds in the others). Each scenario was seen by between 14-16 participants.

For all scenarios, participants were asked “Now, considering each of the possible trophies on the course, please respond on the 1-7 scale: how likely is it that the participant is trying to reach this trophy?” Likert scale options were as follows:

- “1 (Definitely not this trophy)”
- “2 (Very unlikely)”
- “3 (Somewhat unlikely)”
- “4 (Equally likely or unlikely)”
- “5 (Somewhat likely)”
- “6 (Very likely)”
- “7 (Definitely this trophy)”

Note, due to a mismatch, some participants in the `Generic` condition saw one stimuli from `Different-Colored Keys` and vice versa. We reversed fixed this partway through the study.

E. LLM Baselines Experiment

We use the gpt-4-0314 model (i.e., a fixed, non-updatable checkpoint from March 14 for maximal reproducibility) in the Chat Completion API. For each stimuli, we sample until we get 20 well-formed responses, rejecting samples that are not in the correct format.

Example LLM Prompt for World Configuration

Input:

There are three trophies placed throughout this obstacle course: gold, silver, or bronze.
The gold trophy is behind a green door.
The silver trophy is behind a yellow and red door.
The bronze trophy is behind a red door.

Keys to unlock the doors are placed around the obstacle course. On this course, each door can be unlocked by any key as long as it is of a different color than the door.

You see Alice picking up a red key. What would be the likely trophy Alice wants to get?"

output:

```
{"agent" : ["Alice"],
"goals" : ["gold", "silver", "bronze"],
"locations" : {"gold" : "Room A", "silver" : "Room B", "bronze" : "Room C"},
"obstacles" : {"Room A" : ["G"], "Room B" : ["Y", "R"], "Room C" : ["R"] },
"keys" : ["g", "y", "r"],
"max_obstacle" : 2,
"keys_per_door": 1,
"len_key" : 3,
"goal_count" : 3,
"observation_type": "has_objects",
"observation": "(exists (?k - key) (and (has Alice ?k) (iscolor ?k red)))"}
```

Example LLM Prompt for PDDL Domain Definition

In this example we constrain the model such that the LLM can only modify the unlock action of the PDDL domain definition. Our goal is to show that with few-shot prompting in the color-same condition, the LLM can generalize to generate action definition for color-different condition.

¹We exclude 6 participants due to large deviations from average rating

input:

Alice is playing a treasure game to pick up exactly one trophy. There are 3 trophies: gold, silver and bronze, located in room A, room C and room B. There is 1 yellow and 1 green lock for Room A, a red lock for room B and a yellow lock for room C.

Alice needs to unlock the locks to enter the rooms. Each lock can be unlocked with a key with the same color. Each key can only be used once.

You see Alice go over to pick up a red key.

It's unclear what Alice would do next.

output:

```
(:action unlock
:parameters (?a - agent ?k - key ?d - door)
:precondition (and (has ?a ?k) (locked ?d)
  (exists (?c - color) (and (iscolor ?k ?c) (iscolor ?d ?c)))
  (or (and (= (xloc ?a) (xloc ?d)) (= (- (yloc ?a) 1) (yloc ?d)))
    (and (= (xloc ?a) (xloc ?d)) (= (+ (yloc ?a) 1) (yloc ?d)))
    (and (= (- (xloc ?a) 1) (xloc ?d)) (= (yloc ?a) (yloc ?d)))
    (and (= (+ (xloc ?a) 1) (xloc ?d)) (= (yloc ?a) (yloc ?d)))))
:effect (and (not (has ?a ?k)) (not (locked ?d)) )
)
```

F. An example of LLM baseline result

Here we provide an example of GPT-4's response in the baseline condition on the *Color-different* task variant.

input:

There are three trophies placed throughout this obstacle course:
gold, silver, or bronze.

The gold and silver trophies are both located behind a yellow door.
To go to the bronze trophy, you need to go through a red door.

Keys to unlock the doors are placed around the obstacle course.
On this course, each door can be unlocked by any key that is a different color than the door. For example, yellow keys cannot unlock yellow doors, but can unlock doors in green or red.

In this game, players lose points if they pick up more keys than needed.
You may assume that players will use all the keys they pick up.

Alice is a participant on this course. You see Alice pick up a yellow key.

It's unclear what Alice would do next.

How likely is Alice picking up the trophies? Please rate from 1 to 7
and explain.

output:

Gold Trophy: 6
Silver Trophy: 6
Bronze Trophy: 1

Explanation:

Since Alice has picked up a yellow key, she can now unlock either the gold or silver trophy, since both are located behind a yellow door. This makes it quite likely that she will pick up one of those two trophies (rating them both at 6).

However, picking up a yellow key means she cannot unlock the red door, which leads to the bronze trophy. Therefore, it is very unlikely that Alice will pick up the bronze trophy (rating it 1).

It is clear from this example that GPT-4 shows unsatisfactory understanding of the setup of the *Color-different* variant. It is stated that yellow keys cannot unlock yellow doors, so it is unlikely that Alice would pick up the gold or silver trophy, but GPT-4 still judges these two possibilities to be highly likely. On the other hand, both humans and our NIPE model assign significantly more likelihood to the bronze trophy (and roughly equal low likelihoods to the gold and silver trophies) in this example.

G. Sample stimuli for each experimental condition

Generic:

There are three trophies placed throughout this obstacle course: gold, silver, or bronze.
To get to the gold trophy, you need to unlock 2 different doors. To get to the silver trophy, you need to unlock 1 door. The bronze trophy isn't behind a door.
Keys to unlock the doors are placed around the obstacle course. On this course, each door has exactly 1 lock on it. All of the keys on this course work on all of the locks. However, each key can only be used one time.
Alice is a participant on this course. You see Alice go over and pick up two keys.

Color-same:

There are three trophies placed throughout this obstacle course: gold, silver, or bronze.
The gold and silver trophy trophies are both located behind a yellow door. To get to the yellow door, you first need to go through a green door. To go to the bronze trophy, you need to go through a red door.
Keys to unlock the doors are placed around the obstacle course. On this course, each door is unlocked by a key of the corresponding color.
In this game, players lose points if they pick up more keys than needed. You may assume that players will use all the keys they pick up.
Alice is a participant on this course. You see Alice pick up a yellow key.

Color-different:

There are three trophies placed throughout this obstacle course: gold, silver, or bronze, which are placed in rooms labeled A, B, and C, respectively.
Room A is behind a door with a yellow and green lock on it. Room B is behind a door with a red lock on it. Room C is behind a door with its own yellow lock on it.
Keys to unlock the doors are placed around the obstacle course. On this course, each door can be unlocked by any key that is a different color than the door. For example, yellow keys cannot unlock yellow doors, but can unlock green and red doors
In this game, players lose points if they pick up more keys than needed. You may assume that players will use all the keys they pick up.
Alice is a participant on this course. You see Alice pick up a yellow key.
It's unclear what Alice would do next.

Spatial:

There are three trophies placed throughout this obstacle course: gold, silver, or bronze. Alice is a participant on this course.

The gold trophy is currently 1 step East of Alice. The silver trophy is 3 steps South from her. The bronze trophy is 3 steps East from the silver trophy.
You see Alice walking 3 steps East.