| Domain | Humans $n$=5 | RL Agents $\alpha$=50 | Persistence ($r$) | | | Persistence ($q$) | | | Search Temp. ($T$) | | | Heuristic ($h$) | | | |
| | | | 1 | 2* | 4 | 0.8 | 0.9 | 0.95* | 2 | 10* | 50 | Manhat.* | Maze.Dist. | Goal.Count. | $h_{\text{add}}$* |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Doors, Keys, Gems | 0.79 | 0.58 | 0.60 | 0.73 | 0.73 | 0.53 | 0.60 | 0.73 | 0.67 | 0.73 | 0.77 | 0.73 | 0.90 | – | – |
| Block Words | 0.73 | 0.82 | 0.90 | 0.87 | 0.90 | 0.70 | 0.83 | 0.87 | 0.83 | 0.87 | 0.87 | – | – | 0.43 | 0.87 |

**Table 1: Robustness to model mismatch.** Top-1 accuracy of SIPS at the third time quartile (Q3), evaluated on data generated by humans, RL agents, and mismatched models. We ran SIPS assuming $r$=2, $q$=0.95, $T$=10, and a Manhattan ($h_{\text{add}}$) heuristic for Doors, Keys, Gems (Block Words). Matched parameters are starred (*).

1  We thank the reviewers for engaging carefully with our paper, and for providing helpful and constructive feedback. We
2  commit to addressing the minor issues raised and adding the suggested references.

3  **R1, R2, R3, R4 raised concerns about whether our model of boundedly-rational planning is robust to plans**
4  **generated by humans and other models**. We appreciate these concerns, and agree that it best to perform a user study
5  as R2 suggests, and to avoid the circular evaluation pointed out by R3. In response, we have performed a series of
6  robustness experiments (**Table 1**), showing that SIPS is robust to data from 5 pilot human subjects (30 trajectories per
7  subject), a Boltzmann RL agent, mismatched parameters $r$, $q$, $T$, and $h$. While performance can degrade with mismatch,
8  this is partly due to the difficulty of inference from highly random behavior (e.g. $q$=0.8, $h$=Goal Count). In fact, when
9  mismatched parameters are *more* optimal, performance can *improve* (e.g.$h$=Maze Dist.). Importantly, SIPS does well
10  on human data (Top-1=0.78 / 0.73), showing robustness even when the planner is unknown. **We have also conducted**
11  **preliminary experiments showing that human goal inferences correlate highly with SIPS** (Pearson's $r = 0.85$).
12  We will expand on these experiments in the final paper with more domains and cross-method comparisons.

13  **R1 asked about the motivation for using a negative binomial for search budgets** $\eta \sim \text{NB}(r, q)$. To explain, $r$ and
14  $q$ model the persistence of a planner who may give up after expanding each node. 1-$q$ is the probability that the planner
15  considers giving up, while $r$ is the number of times the planner has to consider giving up before actually giving up. This
16  captures the intuition that people are more likely to give up the longer we plan, while still exhibiting some persistence.
17  **R1 also expressed concerns about our relaxation of A* search**, where a node $s$ is sampled from the open list for
18  expansion according to $P_{\text{expand}}(s) \propto \exp(-f(s, g)/T)$. To explain, this captures how agents may fail to choose the
19  best cognitive action (i.e. node expansion) during search, following a Boltzmann distribution. $T$ controls how informed
20  the search is, with $T$=0 (least informed) equivalent to breadth-first search, $T$=$\infty$ (most informed) equivalent to standard
21  A*. Search is sound, because all nodes are eventually expanded, and path costs updated. Similar modifications are
22  made in both [**1**] and [**2**]. We will add these references in the final paper.

23  **R2 and R3 were also concerned whether our modeling assumptions are overly specific** *(R2. Q1, R3. Weaknesses)*.
24  We agree that there are many types of planning that go beyond the scope of this paper. As we note in Future Work,
25  this motivates extensions to domains like task and motion planning. To R3's point on optimal RL actions without
26  explicit planning, we agree this is a good model when humans are well-practiced. Still, when inferring plans for novel
27  tasks with sparse rewards, we believe explicit planning is a better model, and will gladly add discussion on this point.
28  **Planners can also exhibit many other bounds, including limited memory (R2) or inference capacity (R3).** We
29  see modeling such bounds as important future work. While our current model no doubt approximates human planning
30  by leaving out these bounds, humans are likely to use similarly approximate models when inferring the goals of others.
31  As **Table 1** shows, this approximation is enough to ensure reasonable robustness in the domains we consider.

32  **R1 asked about the importance of probabilistic programming for our research.** While R1 is right that the basic
33  particle filtering approach could be implemented manually with some additional work, our use of custom proposal
34  programs [**3**] and involutive MCMC [**4**] in the SIPS rejuvenation moves requires computing importance weight ratios for
35  all random choices in both the proposals and the model. Without Gen's automated weight computation, implementing
36  this would be highly tedious and error prone, akin to implementing back-propagation for a VAE without TensorFlow or
37  PyTorch. A manual implementation would also be much more difficult to extend to complex goal priors and planners.
38  We are happy to emphasize this in the final paper.

39  **R2 asked for clarification about whether our approach is tied to deterministic settings.** While all the evaluated
40  domains were deterministic, our framework also supports Probabilistic PDDL [**5**], and can be readily extended to
41  stochastic domains, e.g. by determinization as R2 suggests. Our use of a replanning model also means that it is not
42  as strongly tied to deterministic environments as R3 contends: if the agent encounters a state that it did not plan for
43  (e.g. randomly failing to pick up a block), it simply replans from that new state, which we believe to be a reasonable
44  approximation of human behavior in mostly deterministic domains. We will clarify these points in the final paper.

45  **R4 asked why unbiased BIRL performs better in the Taxi environment, and why oracle BIRL sometimes per-**
46  **forms better for top-1 accuracy.** This is because SIPS can suffer from particle collapse, whereas BIRL can perform
47  exact inference if it has good Q-value estimates (true for both oracle BIRL and the Taxi domain's small state space).
48  SIPS can be improved by increasing the particle count, using rejuvenation, or variance control techniques. We have
49  since implemented some of these techniques (e.g. residual resampling), and will update our results accordingly. **R4**
50  **also asked how well SIPS performs when the agent does not replan.** This is addressed in the supplement, where we
51  show that SIPS performs reasonably well on trajectories from an optimal full-horizon planner.