# Inferring Rewards from Language in Context

**Jessy Lin**◇    **Daniel Fried**♣    **Dan Klein**◇    **Anca Dragan**◇

◇ University of California, Berkeley

♣ Carnegie Mellon University

{jessy_lin, klein, anca}@berkeley.edu, dfried@andrew.cmu.edu

## Abstract

In classic instruction following, language like "I'd like the JetBlue flight" maps to actions (e.g., selecting that flight). However, language also conveys information about a user's underlying reward function (e.g., a general preference for JetBlue), which can allow a model to carry out desirable actions in new contexts. We present a model that infers rewards from language pragmatically: reasoning about how speakers choose utterances not only to elicit desired actions, but also to reveal information about their preferences. On a new interactive flight–booking task with natural language, our model more accurately infers rewards and predicts optimal actions in unseen environments, in comparison to past work that first maps language to actions (instruction following) and then maps actions to rewards (inverse reinforcement learning).

## 1 Introduction

Language is a natural interface for systems like robots or personal assistants that interact with human users. One way to interpret language in these interactive settings is to train an instruction following agent: a model that learns to map commands like "go three steps forward to the door" to a sequence of actions in context (e.g., Branavan et al. 2009; Tellex et al. 2011, *inter alia*). Instructions describe *how* an agent should act in an immediate context, but to build models that can generalize—carrying out a user's goals in new contexts and learning user preferences over repeated interactions—agents should also infer *why* actions are taken. Grounding language to *reward functions* extends the standard instruction following setup in this way, representing the goals and preferences that underlie actions, and allowing agents to autonomously carry out correct actions in new contexts (e.g., Fu et al. 2019).

However, when people interact with systems they often primarily aim to achieve specific tasks,
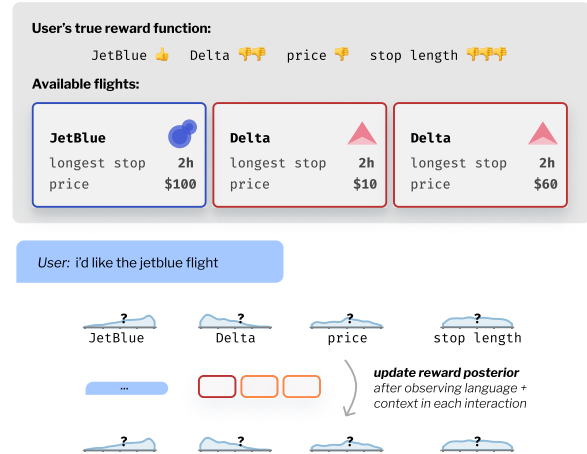


Figure 1: **When people instruct agents with language like "I'd like the JetBlue flight," both their desired actions and the language itself reveal information about rewards.** From the referenced flight itself, a model would guess that the user may prefer expensive JetBlue flights. Reasoning jointly with language reveals that JetBlue is the more salient preference, and the model should still have uncertainty about whether expensive flights are generally preferred. JetBlue may have been more important than a preference for cheap flights, but the user may still prefer cheap flights, all else equal. Over repeated interactions with the user in new contexts, the model can continually refine its estimates of the user's preferences.

rather than literally describing their preferences in full. How do we infer general goals and preferences from utterances in these settings? Consider a flight booking agent like the one in Figure 1. By inferring the user's reward function (indicating their preference for carrier, price, and other flight features) beyond just selecting the right flight, such a system would be able to autonomously book flights on behalf of the user in other instances. To do so, the system might use the actions the user commands as evidence about what they prefer, recovering rewards from actions using (language-free) techniques like inverse reinforcement learning (IRL; Ng and Russell 2000). For example, the system can select a flight the user might like in a new instance by

matching features from their past flight bookings.

The key idea of our work is that the *way* that a user refers to their desired actions with language also reveals important information about their reward: the fact that they said "the JetBlue flight" and not "the expensive flight" conveys what matters to them. Intuitively, in settings with repeated interactions, utterances are optimized to communicate information that is generalizable—implicitly helping listeners make useful inferences for acting on a longer horizon. We implement this idea with a pragmatic model of how speakers (humans) generate such language: speakers choose utterances that both elicit reward-maximizing actions in a particular context and faithfully describe the reward. Given an utterance, our model infers that the most likely rewards are the ones that would have made a speaker likely to choose that utterance.

To evaluate our model, we construct and release a dataset for mapping language to rewards, FLIGHT-PREF, containing natural language utterances from humans with underlying preferences. Humans interact in a multi-turn flight booking game similar to Figure 1, where we provide a "user" player with a reward function representing flight preferences. The goal of the game is for the user to communicate these preferences in natural language to an "assistant" player, who is tasked with booking preferred flights for the user. We present this dataset as a challenging benchmark for reward learning from language and interaction.

In our experiments, we show that our model can infer reward functions from natural language, improve reward estimates consistently over repeated interactions, and use inferred rewards to accurately select optimal actions in held-out environments. Our full model obtains relative accuracy improvements of 12% when compared to models that only treat language as descriptions of actions.[1]

## 2 Related Work

**Instruction following.** A long line of work on grounded instruction following has developed various methods for producing actions from language, including approaches that use intermediary structured semantic representations (MacMahon et al., 2006; Tellex et al., 2011; Chen and Mooney, 2011; Matuszek et al., 2013; Artzi and Zettlemoyer, 2013; She et al., 2014; Thomason et al., 2015; Wang et al.,

2016; Fried et al., 2018a; Arumugam et al., 2017; Suhr et al., 2018) or map directly to primitive actions (Branavan et al., 2009; Andreas and Klein, 2015; Mei et al., 2016; Bisk et al., 2016; Misra et al., 2017; Guu et al., 2017; Suhr and Artzi, 2018; Anderson et al., 2018; Shridhar et al., 2020). All of these approaches interpret any given utterance (instruction) solely in the context that elicited the utterance, producing one particular sequence of actions. The method we present extends these approaches, using utterances to infer the rewards that underlie the actions that should be taken across a range of environments: both the context that elicited the utterance, and other unseen environments.

**Reward learning.** The majority of work on reward learning has been in the robotics and reinforcement learning communities and has not incorporated language, rather using techniques such as inverse reinforcement learning (IRL; Ng and Russell 2000; Ratliff et al. 2006; Ziebart et al. 2008; Hadfield-Menell et al. 2017; Jeon et al. 2020) to infer the rewards that underlie human demonstrations of actions. Even works that incorporate language into reward learning also take this primarily action-centric approach: either by using datasets pairing utterances with *trajectories* and using (language-free) IRL to then recover reward functions from trajectories (MacGlashan et al., 2015; Fu et al., 2019), or learning an instruction-following model guided by a language-conditioned discriminator (Bahdanau et al., 2019). The language in these settings are unambiguous commands, giving a complete description of a goal (e.g., "go to the red door"). In contrast, we are concerned with language used to guide agents in repeated interactions (where language may be a partial or ambiguous mix of instructions and reward descriptions).

**Pragmatics.** A long line of work on pragmatics (Grice, 1975), particularly in the Rational Speech Acts (RSA) framework (Goodman and Frank, 2016), has developed computational models for inferring the behavior or belief that a speaker wishes to induce in a listener. However, the majority of this work has only focused on single-turn interactions, where an utterance conveys an action in a single context, e.g., choosing the correct referent in signaling games (Golland et al., 2010; Frank and Goodman, 2012; Degen et al., 2013; Monroe et al., 2017; McDowell and Goodman, 2019), interpreting implicatures (Goodman and Stuhlmüller,

---

[1] We release our code and dataset at https://github.com/jlin816/rewards-from-language.
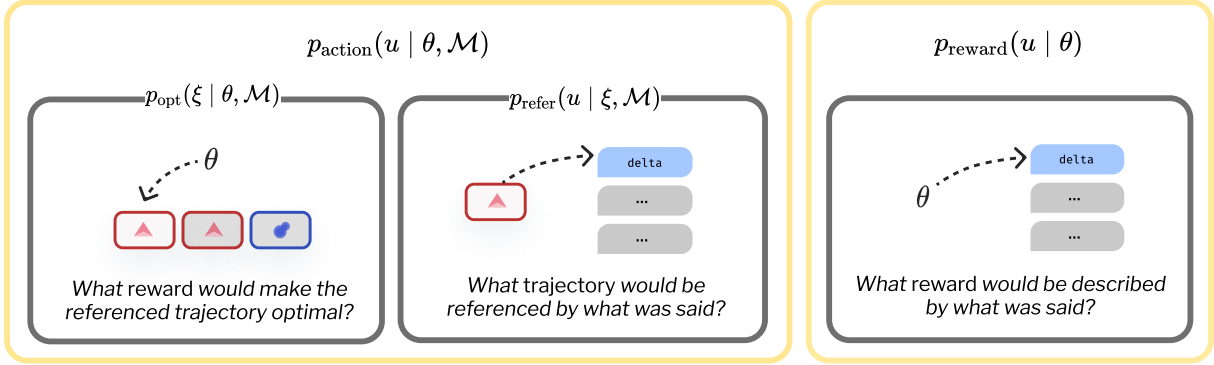
Figure 2: Our model infers rewards by reasoning about how the speaker chose the observed utterance: both to elicit correct actions ($p_\text{action}$) and to describe their reward ($p_\text{reward}$). We illustrate this on the flight domain, where trajectories are a choice of a single flight.

2013; Bergen et al., 2016), or generating (Fried et al., 2018a; Sumers et al., 2021) or interpreting grounded instructions (Fried et al., 2018b). Our work extends this past work by showing that in repeated interactions, listeners can also benefit by reasoning pragmatically about how speakers communicate information about and over longer time horizons.

## 3 Reward Inference from Language

**Problem Formulation.** We parameterize the user's preference as a reward function $r_\theta$ with parameters $\theta$. In our flight booking domain from Figure 1, $\theta$ is a weight vector which specifies preferences over flight features (carrier, price, etc.). We formalize the general reward inference problem as sequence of Markov decision processes (MDPs) $\mathcal{M}_1, \ldots, \mathcal{M}_I$ that share the same reward function $r_\theta$. In each MDP $\mathcal{M}_i$, the agent receives an utterance $u_i$ from the user and must execute a trajectory $\xi$. The agent's goal is to infer $\theta$ over the sequence of interactions, which should allow the agent to execute trajectories with high reward in as-yet unseen contexts.

The agent maintains an estimate over $\theta$ over the course of interactions. We introduce a model $p(\theta \mid u, \mathcal{M})$ that the agent will use to perform Bayesian updates of a posterior over $\theta$:

$$p(\theta \mid u_{1:i}, \mathcal{M}_{1:i}) \propto p(\theta \mid u_i, \mathcal{M}_i)$$
$$\times p(\theta \mid u_{1:i-1}, \mathcal{M}_{1:i-1})$$

In the flight domain, we specialize this formulation to study a one-step MDP (contextual bandit). Trajectories $\xi$ consist of a single action, choosing one of the available flights. Over a series of these rounds where the agent books a flight given the user's utterance $u_i$, the agent must infer the user's

flight preferences $\theta$ to book flights from other unseen sets of options, without explicit language instruction from the user.

### 3.1 Model

Our model, summarized in Figure 2, defines a *rational listener*, $L_2$, which predicts a distribution over rewards $\theta$, conditioned on an utterance $u$ and a context $\mathcal{M}$. (The terminology we use for listeners and speakers follows Bergen et al. 2016.) The rational listener uses Bayesian reasoning about a speaker model, $S_1$, which produces utterances conditioned on a reward function and context:

$$p_{L_2}(\theta \mid u, \mathcal{M}) \propto p_{S_1}(u \mid \theta, \mathcal{M})p(\theta \mid \mathcal{M})$$

Key to our model is that the $S_1$ speaker distribution $p_{S_1}(u \mid \theta, \mathcal{M})$ defines how speakers produce language that functions both to elicit correct actions and describe their underlying reward:

$$p_{S_1}(u \mid \theta, \mathcal{M}) = \alpha p_\text{action}(u \mid \theta, \mathcal{M})$$
$$+ (1 - \alpha)p_\text{reward}(u \mid \theta),$$

where $\alpha$ controls the speaker's "nearsightedness"—how much does the speaker care about the listener choosing the correct action in the *current* context, rather than describing the reward in a context-independent way so that the agent can make good choices in *future* contexts?

**Optimizing for action.** The behavior-optimizing term $p_\text{action}$ specifies that the speaker chooses utterances that elicit reward-maximizing behavior from a listener in the current environment:

$$p_\text{action}(u \mid \theta, \mathcal{M})$$
$$= \sum_\xi p_\text{refer}(u \mid \xi, \mathcal{M})p_\text{opt}(\xi \mid \theta, \mathcal{M}),$$

where the *optimality model* $p_{\text{opt}}(\xi \mid \theta, \mathcal{M})$ specifies the probability the speaker refers to trajectory $\xi$ if their true reward is $\theta$. We can formulate the optimality model with the Boltzmann distribution common in IRL, where speakers are noisily-rational about which trajectories to refer to: $p_{\text{opt}}(\xi \mid \theta, \mathcal{M}) \propto \exp(\beta r_\theta(\xi; \mathcal{M}))$, with rationality parameter $\beta$. This term specifies that utterances are more likely to refer to trajectories that have high reward according to the speaker's $\theta$, compared to other trajectories in $\mathcal{M}$.

Then, for a particular trajectory $\xi$, $p_{\text{refer}}(u \mid \xi, \mathcal{M})$ specifies what utterances are likely to refer to that trajectory. In particular, we model that speakers choose utterances that would make a listener execute that trajectory:

$$p_{\text{refer}}(u \mid \xi, \mathcal{M}) \propto p_{L_{\text{base}}}(\xi \mid u, \mathcal{M})$$

using a *base listener* model $L_{\text{base}}$ of the type common in past work on instruction following. We provide details on $L_{\text{base}}$ in Section 5.

**Optimizing for reward descriptiveness.** Finally, we model $p_{\text{reward}}(u \mid \theta)$, the second term in $P_{S_1}$, with a *base speaker* model, $S_{\text{base}}$, that maps rewards to reward descriptive utterances: $p_{S_{\text{base}}}(u \mid \theta)$[2]. We also provide details on $S_{\text{base}}$ in Section 5.

### 3.2 A Generative Model of Utterances

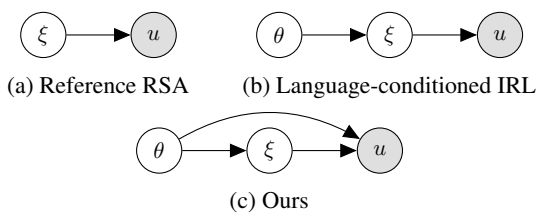(a) Reference RSA     (b) Language-conditioned IRL

(c) Ours

Figure 3: Graphical models contrasting prior work with our model, which models how language utterances $u$ convey both explicit information about the reward $\theta$ and implicit evidence about $\theta$ through the actions they suggest (via trajectories $\xi$). Dependence on $\mathcal{M}$ not shown for visual clarity.

Our account of pragmatic generation can also be viewed as the graphical model in Figure 3(c), where, importantly, the reward influences the utterance both directly and via the action that the

speaker refers to. We define $p(u \mid \xi, \theta, \mathcal{M})$ to be:

$$p(u \mid \xi, \theta, \mathcal{M}) = \alpha p(u \mid \xi, \mathcal{M})$$
$$+ (1 - \alpha) p(u \mid \theta, \mathcal{M})$$

and assume that utterances are reward-descriptive in a way that is independent of the current context, $p(u \mid \theta, \mathcal{M}) = p(u \mid \theta)$.

We can confirm this leads us back to $p_{S_1}$ by marginalizing out $\xi$:

$$p(u \mid \theta, \mathcal{M}) = \sum_\xi p(u \mid \xi, \theta, \mathcal{M}) p(\xi \mid \theta, \mathcal{M})$$
$$= \alpha \sum_\xi \Big( p(u \mid \xi, \mathcal{M}) p(\xi \mid \theta, \mathcal{M}) \Big)$$
$$+ (1 - \alpha) p(u \mid \theta, \mathcal{M})$$
$$= \alpha p_{\text{action}}(u \mid \theta, \mathcal{M}) + (1 - \alpha) p_{\text{reward}}(u \mid \theta)$$

Using this graphical model, we illustrate how our model differs from prior work in similar settings:

**Classic reference game pragmatics collapses belief and behavior.** In general, RSA allows the speaker to optimize for any "utility function," and in the simplest form the utility function optimizes for the listener's belief over world states (Goodman and Frank, 2016). However, in most work on RSA the only relevant world-state belief is belief about behavior, e.g., the referent that should be selected (Figure 3a). Instead, our setting disentangles communication about intended referents in a single context and communication about (reward) beliefs, which influence behavior on longer horizons. Andreas et al. (2017); Sumers et al. (2021) have made the same observation: reference games conflate whether the speaker's objective is to influence beliefs or actions, and modeling the speaker as one or the other produces distinct interpretations of utterances (e.g., speakers that only optimize for correct behavior may do so at the cost of being truthful about the reward).

**IRL assumes all information about the reward function is modulated by the trajectory.** Prior work (MacGlashan et al., 2015; Fu et al., 2019) uses IRL to recover rewards from *trajectories* (e.g., from datasets pairing utterances with trajectories), and then supervising a model with these induced (utterance, reward) pairs. While prior work has not specifically considered pragmatics (i.e., speaker models), their implicit speaker model amounts to assuming that all information about the reward

---

[2]In principle, $p_{\text{reward}}$ could also do pragmatic reasoning to optimize a listener's reward belief, but we did not find an improvement from doing so empirically.

comes from trajectories, as in Figure 3b. In our experiments we compare against a pragmatic version of this action-centric speaker, which is equivalent to setting $\alpha = 1$ in our model (only using $p_{\text{action}}$). In realistic settings where utterances are *not* unambiguous commands like "go to the red door," it becomes important to model how actions and utterances reveal *complementary* information about rewards.

## 4 The FLIGHTPREF Task

We design FLIGHTPREF, a task for reward inference from natural language in the flight booking domain. FLIGHTPREF is designed to simulate a simplified interaction with a flight booking agent, where users communicate with the agent via language to book flights from a set of options. Effective agents must not only learn to book the preferred flight given an instruction in the immediate context (instruction following), but also learn the user's preferences over repeated interactions to book preferred flights in unseen contexts.

We collect a dataset of natural language in a multi-turn game between a user (the "speaker") and an assistant (the "listener" agent). Each flight is represented by a feature vector $\phi(\xi) \in \mathbb{R}^8$ (e.g., features of carrier, price, etc.). We assume the user has a linear reward function with parameters $\theta \in \mathbb{R}^8$, specifying a reward for a particular flight $r_\theta(\xi) = \theta^\mathsf{T}\phi(\xi)$.

In the first round of the game, the user and assistant observe a set of three flight options and the user provides an utterance to describe the flight they want (the optimal flight under the reward function), e.g., "the flight with the most stops." In each of the subsequent rounds, the user and assistant are presented with a new set of three flights. The assistant can either *choose* by guessing the user's preferred flight (under the same reward function), or *prompt* the user for another utterance describing the desired flight in the new set. If the assistant chooses but does so incorrectly, the user is prompted for another utterance describing the correct flight. Both players are penalized if the assistant chooses incorrectly, and earn points if the assistant chooses correctly (with more points for each round the assistant can do so without asking for help). The user is thus incentivized to provide utterances that inform the agent which flight to choose, while enabling long-term success over later rounds.

> one stop that is short
> american is the flight that i want. but i need the flight that is the cheapest and has less stops.
> anything but american
> jetblue one
> i need a flight with any airline but jet blue, price and number of stops are a bad factor for me also. i prefer delta if affordable and low layovers. can you help me?
> even american is undesirable, paying more is important
> i like the flight that is $64

Figure 4: Sample text from the task, exhibiting a diversity of instructive and reward-descriptive language.

### 4.1 Data collection

To collect data for the task, we recruit Amazon Mechanical Turk workers and randomly pair them to play six games (i.e., six different reward functions) of six rounds each. Each game thus consists of 1-6 utterances describing options for the same reward function in different contexts. One person plays the role of the user and the other acts as the assistant. The user has access to a hidden reward function, which is a discretized, randomly-sampled vector $\theta \in \{-1, -0.5, 0, 0.5, 1\}^8$. In total, we collected 2,568 utterances across 813 games, of which we split off the 91 games with the highest score (where the speaker and listener were able to communicate most effectively) for the evaluation set. More details about the data collection process can be found in Section A of the appendix.

A sampling of text is shown in Figure 4. Utterances exhibit a range of phenomena: some users lean towards describing very option-specific features (e.g. "i like the flight that is $64"). Other users attempt to describe as much of their reward function as possible (e.g. "i need a flight with any airline but jetblue,...")—we note that even when they did so, the user's tradeoffs between features remain ambiguous. Many of the utterances are neither fully option-specific nor fully reward-descriptive: instructions like "one stop that is short" both instruct the agent which flight to select in the present context, while communicating some generalizable (but incomplete) information about the user's preferences.

## 5 Model Implementation

Our pragmatic model (Section 3.1) relies on base listener and speaker models $L_{\text{base}}$ and $S_{\text{base}}$. In this section, we describe implementations of these models for the FLIGHTPREF dataset. To train the base models, we use the speaker-side data of (utterance,

option set, reward function) tuples from each round. Our base listener and speaker models assume that the utterances are generated conditionally independently given the reward; we capture the dynamics of multiple turns in the posterior reward inference. Both base models learn neural encodings of utterances $u$, actions $\xi$, and rewards $\theta$, and produce distributions by applying softmax functions to inner products between these encodings. We use $\xi^*$ to denote the optimal action in each context, i.e., $\xi^* = \arg\max_\xi r_\theta(\xi)$.

**Base listener model.** The base listener model $L_{\text{base}}$ is defined using inner product similarities between learned representations of actions $\xi$ produced by an MLP encoder, and learned representations of utterances produced by a BERT-base (Devlin et al., 2019) encoder:

$$p_{L_{\text{base}}}(\xi \mid u, \mathcal{M}) \propto \exp(\text{MLP}_{L_{\text{base}}}(\xi) \cdot \text{BERT}_L(u))$$

where the distribution is normalized over all actions (flights) available in the context, $\xi' \in \mathcal{M}$.

We set the rationality parameter $\beta = \infty$ in $p_{\text{opt}}$ as speakers tend to refer primarily to the optimal option in our domain.

**Base speaker model.** The base reward speaker model $S_{\text{base}}$ is defined using an inner product between representations of rewards $\theta$ from an MLP encoder, and utterance representations from a BERT encoder:

$$p_{S_{\text{base}}}(u \mid \theta) \propto \exp(\text{MLP}_{S_{\text{base}}}(\theta) \cdot \text{BERT}_S(u)/\tau)$$

where $p_{S_{\text{base}}}$ is normalized over a set of utterances taken from the training data (see Section C in the appendix), and $\tau = 3$ is a temperature parameter.

**Training.** We fine-tune all model parameters, including the parameters of the initially-pretrained BERT utterance encoders in the listener and speaker on $(u, \xi, \mathcal{M})$ pairs from the training data using the AdamW optimizer (Kingma and Ba, 2015; Loshchilov and Hutter, 2019). The listener and speaker models are trained separately, without sharing any parameters between the encoders used in the two models. We independently train 5 random seeds of each base model and ensemble them together in evaluation by averaging their output probabilities, which we found improved performance of all models (both our full model and baselines). See Section C in the appendix for details and model hyperparameters.

| Method | Held-out accuracy (%) |
|---|---|
| *Oracle models* (infer $k$ features perfectly) | |
| $k = 1$ | 43.0 |
| $k = 2$ | 51.5 |
| $k = 3$ | 60.2 |
| $k = 4$ | 64.7 |
| Action-only | $52.8 \pm 0.97$ |
| Reward-only | $57.8 \pm 0.95$ |
| Action + reward (Ours) | $59.1 \pm 0.96$ |

Table 1: Average held-out accuracy averaged over all evaluation rounds, with standard error of the mean indicated. Our full action+reward model significantly outperforms action-only and reward-only models (with $p < .05$ using the paired bootstrap test). Held-out accuracy is also shown for oracle models that infer $k$ (randomly-chosen) features of the reward perfectly and maintain a uniform distribution over the other features.

**Pragmatic inference** We follow previous work (Fried et al., 2018a; Monroe et al., 2017) and approximate the $S_1$ distribution by normalizing over a fixed set of utterances: the de-duplicated set of short utterances (less than 8 tokens, making up the majority of utterances) with no digits from the training data. We implement the full pragmatic model $p_{L_2}(\theta \mid u, \mathcal{M})$ in Pyro (Bingham et al., 2018) and use importance sampling to generate samples from the posterior over rewards. Given our dataset collection procedure (where we uniformly sample rewards), we model an uniform prior over rewards $p(\theta \mid \mathcal{M})$ for the first interaction.

## 6 Experiments

We evaluate models in the same repeated turn setup that humans carried out in the task. For each game, models play the role of the listener in that game, updating the reward posterior (Section 3.1) after observing the utterance and option set in each round. Our goal is to estimate rewards that allow the agent to carry out the person's preferences: choosing the optimal option (flight) in unseen contexts (sets of flight options). To that end, we directly compare models on **held-out accuracy**: on 1,000 randomly-generated sets of three options, how often the model's estimate of the reward, $\hat{\theta}$, selects the option that is optimal under the true reward.[3] We use the model's reward posterior mean as the estimate, $\hat{\theta} = \mathbb{E}_{p_\theta}\theta$. We additionally provide com-

---

[3] Note that when collecting the dataset, we also tested human listeners's ability to generalize, but only had them select an option on a single unseen option set—the next one in the sequence—to make data collection tractable.
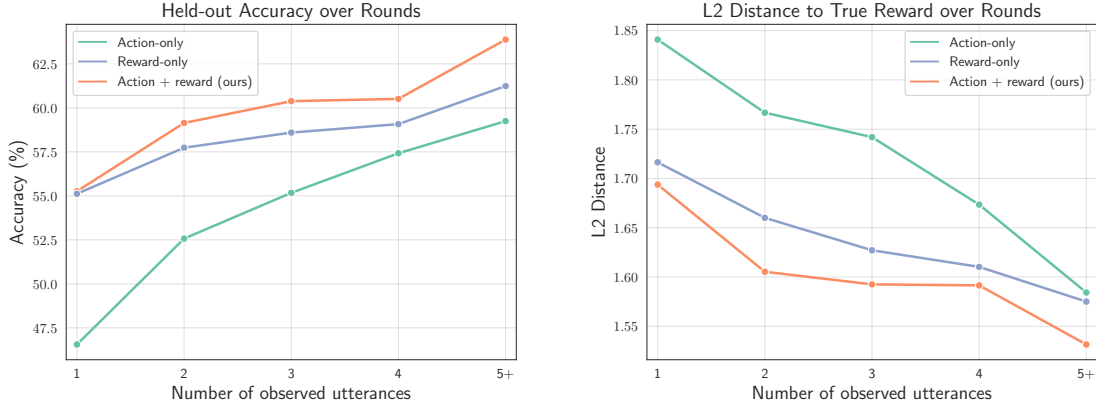
Figure 5: **Multi-turn performance on held-out accuracy (left) and L2 distance to the true reward (right).** We show the performance of each model for varying numbers of observed utterances for a given reward. We combine five- and six-utterance rounds as there were $< 25$ samples in each of these bins. Our full action+belief model substantially outperforms an action-only model at all numbers of utterances ($p < .05$), and performs comparably to or better than a belief-only model, with statistically significant benefits for 5+ utterances ($p < .05$).

parisons of **reward L2 distance** between the estimated reward and the true reward as a context-independent metric: $\sqrt{\sum_{i=1}^{8}(\hat{\theta}_i - \theta_i^*)^2}$, where $\theta^*$ is the true reward.

For our full *action + reward* model, we set the nearsightedness parameter $\alpha = 0.5$ for all posterior updates. We compare to an *action-only* model that uses only $p_{\text{action}}$ (i.e., setting $\alpha = 1.0$). This model is representative of approaches from past work on language-conditioned reward learning (e.g., Mac-Glashan et al. 2015; Fu et al. 2019) that infer rewards purely from the actions that utterances refer to. We also compare to a *reward-only* model that uses only $p_{\text{reward}}$ (inferring rewards purely from the utterance, without conditioning on actions, i.e., setting $\alpha = 0.0$). For comparison to versions of our approach that remove pragmatic modeling, see Section D.1 in the appendix.

### 6.1 Overall Results

In Table 1 we compare all models on held-out accuracy averaged over all rounds in the evaluation set (for each round, having observed all previous rounds in that game). Note that because held-out accuracy is assessed by the proportion of randomly-generated flight sets (out of 1,000) where the true reward function and the inferred reward function pick out the same optimal flight, it is significantly more difficult than achieving high accuracy on a single three-choice instance.

Our full action+reward model achieves a held-out accuracy of 59.1%, +6.3% over the action-only model and +1.3% over the reward-only model, indi-

cating that combining both sources of information allows better inference of rewards that enable optimal actions in novel contexts. For reference, an oracle baseline that infers the value of $k$ randomly chosen features perfectly and is uniform on the other features obtains the following held-out accuracies: $k =$**1** (43%), **2** (51%), **3** (60%), **4** (65%), showing that our model is able to attain similar generalization performance even in the presence of uncertainty (without receiving oracle information about the true value of any feature).

We analyze why our model benefits from both components in Section 6.3, and discuss potential for further improvements in Section 6.4.

### 6.2 Learning over Multiple Interactions

We explore how each model's reward inferences change as more observations are obtained over the course of a game. In Figure 5, we plot held-out accuracy and L2 distance to the true reward as a function of number of observed utterances. Our model outperforms the action-only and reward-only models for all numbers of observed utterances.

**Relying on explicit information from language is most important when there are few observations.** While our full action+reward model improves substantially over the action-only model at all points, this improvement generally decreases as more utterances are observed (Figure 5). Conversely, the improvement of the full model over reward-only generally increases. Qualitatively, we observe that this occurs because utterances tend to mention the most extreme features of the reward

**Southwest**

| | |
|---|---|
| arrival time | 0.0 |
| longest stop | .96 |
| # stops | .50 |
| price | .59 |

**Delta**

| | |
|---|---|
| arrival time | .92 |
| longest stop | .56 |
| # stops | .50 |
| price | .07 |

**Delta**

| | |
|---|---|
| arrival time | .08 |
| longest stop | .80 |
| # stops | 0.0 |
| price | .78 |

*'want a long time before meeting"*

Arrival time

Action-only

Belief-only

Ours

(a) Both the described action (the referenced flight is the one with the highest arrival time) and the explicit reward description in the utterance provide evidence that the user's true reward on arrival time is positive, leading the posterior in our model to (correctly) place more probability mass on positive values of this feature.

**Delta**

| | |
|---|---|
| arrival time | .52 |
| longest stop | .64 |
| # stops | .50 |
| price | .54 |

**JetBlue**

| | |
|---|---|
| arrival time | .24 |
| longest stop | .92 |
| # stops | .50 |
| price | .94 |

**JetBlue**

| | |
|---|---|
| arrival time | 1.0 |
| longest stop | .08 |
| # stops | .25 |
| price | .37 |

*'i sort of like cheaper flights and fewer stops"*

Arrival time     # stops     Price

Action-only

Belief-only

Ours

(b) Evidence from actions and from the utterance complement each other: the action-based model captures that rewards that are positive on arrival time make the selected flight optimal, even though it is unmentioned, while the reward-based model captures evidence about the reward from the user's utterance.

Figure 6: **Real examples showing reward posteriors of each model after observing the given utterance and options.** We sample from the posterior over rewards and visualize the marginal probability distributions for particular features using kernel density estimation. The true reward value for the feature is marked with a red line and the posterior mean for the feature with a blue line.

function, which allow our model to estimate the values of these important features. When there are few observations, inferring reward information from utterances in this way is more informative than using only the option implied by the user's utterance, which does not disambiguate between rewards that select the same option (a commonly discussed problem in IRL; Ziebart et al. (2008)).

**Inferring evidence from actions is most important when there are more observations.** We observe that the action-only model improves more consistently over rounds. Qualitatively, the information that utterances provides about rewards is correlated across multiple rounds—speakers frequently mention salient reward features, whereas actions consistently provide new information about all features. This is particularly pronounced in our domain, due to a relatively small feature and action space. In other more complex domains, actions might provide even more benefits as they provide *fine-grained* information about reward values and tradeoff boundaries that are more difficult to communicate precisely in language.

### 6.3 Analyzing the Benefits of Combining Actions and Rewards

In this section, we investigate *why* our model benefits from both the action and reward models.

**A single utterance and context can provide useful evidence to both models.** In Figure 6, we show the reward posteriors for each model after a single update on a round (starting from a uniform prior). In Figure 6a, we observe how the action- and reward-only models can make correlated updates on an utterance and context where both the action (a flight with a high value on arrival time) and the utterance provide evidence about the arrival time feature. This leads our model's posteriors to aggregate more probability mass on positive values of that feature. In Figure 6b, we show how each model can make inferences about different features for the same context—the action-only model inferring positive values for arrival time given the observed flight and the reward-only model updating on flight price and stops. Our model posterior aggregates information from both.

**Some utterances are primarily "nearsighted," and others primarily "farsighted."** Another reason our full model improves is because some utterances are particularly "farsighted"—mentioning a great deal of explicit information about the reward (which the action-only model cannot take advantage of)—while other utterances are more "nearsighted"—specialized to the particular action, e.g., saying just enough to uniquely identify the optimal flight. Sorting the utterances by difference in accuracy between the action-only and reward-only models confirms that they exhibit qualitatively different phenomena: examples where the reward-only model helps the most are highly reward-

descriptive (e.g., "if i had a choice, i would never fly with delta and american! get me jetblue or southwest...") while examples where the action-only model helps most have less informative utterances (e.g.,"the cheaper the better"). Our full model is able to handle both kinds of language use.

To further analyze the influence of the action and reward component, we evaluate an oracle model that *switches* between the action-only and reward-only models, choosing the model with highest held-out accuracy in each round. This model outperforms our action+reward model (improving from 59.1 to 62.9% on overall held-out accuracy), suggesting that further improvements could be obtained by integrating evidence from the two models. Doing so optimally is challenging in our setting: when a user says "i like the cheap jetblue flight," do they mean to say they like JetBlue generally, or just that they want to choose a desirable flight that happens to be uniquely identified by JetBlue? Future work might explore adaptively switching policies (e.g., using the utterance, or knowledge about the user).

### 6.4 Inference Improves with Known Actions

While our base models have fairly high performance (e.g., the base listener model $L_{base}$ has an average accuracy of 74% at selecting the optimal choice in each option set that has an utterance in the evaluation data), they naturally have some errors which lead to errors in reward inference. We test the influence of this underlying prediction error by skipping posterior updates on all rounds where the base listener predicts the incorrect option for the true reward function. This change improves held-out accuracy by 6% over the reward-only model after six observations (+4% from the original gap), indicating (1) that dataset affords future work on improved instruction following models and (2) that our reward inference procedure benefits from base model improvements.

We note that in our task design, the user does not provide a demonstration (i.e., a choice of flight) to the model. However, if it is convenient to obtain demonstrations from users (e.g., a flight booking interface could let the person click on the flight they want in addition to specifying what they want in natural language), demonstrations would effectively serve as an oracle instruction-following model for that context, which could be incorporated into our full reward inference model.

## 7    Discussion & Conclusion

We presented a method for using natural language to infer reward functions: representing the goals, preferences, and intents underlying action.

Conceptually, our work builds on previous work on language grounding by exploring how language serves a dual purpose. Utterances can refer directly to actions to be taken, as studied in instruction following. Beyond that, they communicate information about "why" those actions should be taken, and what actions may be desirable in new contexts. To build language-guided agents that can interact with people over longer horizons, it may be useful to model this relationship between language, actions, and rewards.

Furthermore, language is *ambiguous* about both actions and goals. Standard settings for studying pragmatics (e.g., reference games) address how to resolve ambiguity about what object or action the speaker is choosing to refer to. We have explored how these settings can be extended by considering the preferences underlying those choices. We introduced FLIGHTPREF, a new dataset of naturalistic interactions between people in a multi-turn flight booking game. FLIGHTPREF uses held-out accuracy as a metric for evaluating interpretation success beyond selecting the right action in a single environment.

Future work can build on the task by 1) learning or evaluating with more complex reward functions (e.g., using deep reward representations); 2) exploring how people communicate about their real preferences and modeling a natural prior (e.g., that people tend to prefer cheaper flights), instead of providing annotators with ground-truth preferences; 3) allowing other ways to handle uncertainty, e.g., leveraging the reward posterior to interactively learn to ask; or 4) extending these approaches to other domains where modeling goals and preferences may be important (e.g., language-conditioned robotics).

### Acknowledgements

# References

Peter Anderson, Qi Wu, Damien Teney, Jake Bruce, Mark Johnson, Niko Sünderhauf, Ian D. Reid, Stephen Gould, and Anton van den Hengel. 2018. Vision-and-language navigation: Interpreting visually-grounded navigation instructions in real environments. In *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018*, pages 3674–3683. IEEE Computer Society.

Jacob Andreas, Anca Dragan, and Dan Klein. 2017. Translating neuralese. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 232–242, Vancouver, Canada. Association for Computational Linguistics.

Jacob Andreas and Dan Klein. 2015. Alignment-based compositional semantics for instruction following. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 1165–1174, Lisbon, Portugal. Association for Computational Linguistics.

Yoav Artzi and Luke Zettlemoyer. 2013. Weakly supervised learning of semantic parsers for mapping instructions to actions. *Transactions of the Association for Computational Linguistics*, 1:49–62.

Dilip Arumugam, Siddharth Karamcheti, Nakul Gopalan, Lawson LS Wong, and Stefanie Tellex. 2017. Accurately and efficiently interpreting human-robot instructions of varying granularities. *arXiv preprint arXiv:1704.06616*.

Dzmitry Bahdanau, Felix Hill, Jan Leike, Edward Hughes, Seyed Arian Hosseini, Pushmeet Kohli, and Edward Grefenstette. 2019. Learning to understand goal specifications by modelling reward. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net.

Leon Bergen, Roger Levy, and Noah Goodman. 2016. Pragmatic reasoning through semantic inference. *Semantics and Pragmatics*, 9.

Eli Bingham, Jonathan P. Chen, Martin Jankowiak, Fritz Obermeyer, Neeraj Pradhan, Theofanis Karaletsos, Rohit Singh, Paul A. Szerlip, Paul Horsfall, and Noah D. Goodman. 2018. Pyro: Deep universal probabilistic programming. *CoRR*, abs/1810.09538.

Yonatan Bisk, Deniz Yuret, and Daniel Marcu. 2016. Natural language communication with robots. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 751–761, San Diego, California. Association for Computational Linguistics.

S.R.K. Branavan, Harr Chen, Luke Zettlemoyer, and Regina Barzilay. 2009. Reinforcement learning for mapping instructions to actions. In *Proceedings of the Joint Conference of the 47th Annual Meeting of the ACL and the 4th International Joint Conference on Natural Language Processing of the AFNLP*, pages 82–90, Suntec, Singapore. Association for Computational Linguistics.

David L. Chen and Raymond J. Mooney. 2011. Learning to interpret natural language navigation instructions from observations. In *Proceedings of the Twenty-Fifth AAAI Conference on Artificial Intelligence, AAAI 2011, San Francisco, California, USA, August 7-11, 2011*. AAAI Press.

Judith Degen, Michael Franke, and Gerhard Jager. 2013. Cost-based pragmatic inference about referential expressions. In *Proceedings of the annual meeting of the cognitive science society*, volume 35.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.

Michael C. Frank and Noah D. Goodman. 2012. Predicting pragmatic reasoning in language games. *Science*.

Daniel Fried, Jacob Andreas, and Dan Klein. 2018a. 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)*, pages 1951–1963, New Orleans, Louisiana. Association for Computational Linguistics.

Daniel Fried, Ronghang Hu, Volkan Cirik, Anna Rohrbach, Jacob Andreas, Louis-Philippe Morency, Taylor Berg-Kirkpatrick, Kate Saenko, Dan Klein, and Trevor Darrell. 2018b. Speaker-follower models for vision-and-language navigation. In *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*, pages 3318–3329.

Justin Fu, Anoop Korattikara, Sergey Levine, and Sergio Guadarrama. 2019. From language to goals: Inverse reinforcement learning for vision-based instruction following. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net.

Dave Golland, Percy Liang, and Dan Klein. 2010. A game-theoretic approach to generating spatial descriptions. In *Proceedings of the 2010 Conference on Empirical Methods in Natural Language Processing*, pages 410–419, Cambridge, MA. Association for Computational Linguistics.

Noah D. Goodman and Michael C. Frank. 2016. Pragmatic language interpretation as probabilistic inference. *Trends in Cognitive Sciences*, 20(11).

Noah D Goodman and Andreas Stuhlmüller. 2013. Knowledge and implicature: Modeling language understanding as social cognition. *Topics in cognitive science*, 5(1):173–184.

H. Paul Grice. 1975. Logic and conversation. In Peter Cole and Jerry L. Morgan, editors, *Speech Acts*, volume 3 of *Syntax and Semantics*.

Kelvin Guu, Panupong Pasupat, Evan Liu, and Percy Liang. 2017. From language to programs: Bridging reinforcement learning and maximum marginal likelihood. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1051–1062, Vancouver, Canada. Association for Computational Linguistics.

Dylan Hadfield-Menell, Smitha Milli, Pieter Abbeel, Stuart J. Russell, and Anca D. Dragan. 2017. Inverse reward design. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pages 6765–6774.

Hong Jun Jeon, Smitha Milli, and Anca D. Dragan. 2020. Reward-rational (implicit) choice: A unifying formalism for reward learning. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.

Diederik P. Kingma and Jimmy Ba. 2015. Adam: A method for stochastic optimization. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*.

Ilya Loshchilov and Frank Hutter. 2019. Decoupled weight decay regularization. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net.

James MacGlashan, Monica Babes-Vroman, Marie desJardins, Michael L. Littman, Smaranda Muresan, S. Squire, Stefanie Tellex, Dilip Arumugam, and Lei Yang. 2015. Grounding english commands to reward functions. In *Robotics: Science and Systems*.

Matt MacMahon, Brian Stankiewicz, and Benjamin Kuipers. 2006. Walk the talk: Connecting language, knowledge, and action in route instructions. In *Proceedings of the Twenty-First AAAI Conference on Artificial Intelligence, AAAI 2006, Boston, Massachusetts, USA, July 16-20, 2006*.

Cynthia Matuszek, Evan Herbst, Luke Zettlemoyer, and Dieter Fox. 2013. Learning to parse natural language commands to a robot control system. In *Experimental robotics*, pages 403–415. Springer.

Bill McDowell and Noah Goodman. 2019. Learning from omission. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 619–628, Florence, Italy. Association for Computational Linguistics.

Hongyuan Mei, Mohit Bansal, and Matthew R. Walter. 2016. Listen, attend, and walk: Neural mapping of navigational instructions to action sequences. In *Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence, February 12-17, 2016, Phoenix, Arizona, USA*, pages 2772–2778. AAAI Press.

Dipendra Misra, John Langford, and Yoav Artzi. 2017. Mapping instructions and visual observations to actions with reinforcement learning. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 1004–1015, Copenhagen, Denmark. Association for Computational Linguistics.

Will Monroe, Robert X.D. Hawkins, Noah D. Goodman, and Christopher Potts. 2017. Colors in context: A pragmatic neural model for grounded language understanding. *Transactions of the Association for Computational Linguistics*, 5:325–338.

Andrew Y. Ng and Stuart J. Russell. 2000. Algorithms for inverse reinforcement learning. In *Proceedings of the Seventeenth International Conference on Machine Learning (ICML 2000), Stanford University, Stanford, CA, USA, June 29 - July 2, 2000*, pages 663–670. Morgan Kaufmann.

Nathan D. Ratliff, J. Andrew Bagnell, and Martin Zinkevich. 2006. Maximum margin planning. In *Machine Learning, Proceedings of the Twenty-Third International Conference (ICML 2006), Pittsburgh, Pennsylvania, USA, June 25-29, 2006*, volume 148 of *ACM International Conference Proceeding Series*, pages 729–736. ACM.

Lanbo She, Shaohua Yang, Yu Cheng, Yunyi Jia, Joyce Chai, and Ning Xi. 2014. Back to the blocks world: Learning new actions through situated human-robot dialogue. In *Proceedings of the 15th Annual Meeting of the Special Interest Group on Discourse and Dialogue (SIGDIAL)*, pages 89–97, Philadelphia, PA, U.S.A. Association for Computational Linguistics.

Mohit Shridhar, Jesse Thomason, Daniel Gordon, Yonatan Bisk, Winson Han, Roozbeh Mottaghi, Luke Zettlemoyer, and Dieter Fox. 2020. ALFRED: A benchmark for interpreting grounded instructions for everyday tasks. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June 13-19, 2020*, pages 10737–10746. IEEE.

Alane Suhr and Yoav Artzi. 2018. Situated mapping of sequential instructions to actions with single-step reward observation. In *Proceedings of the 56th Annual Meeting of the Association for Computational*

*Linguistics (Volume 1: Long Papers)*, pages 2072–2082, Melbourne, Australia. Association for Computational Linguistics.

Alane Suhr, Srinivasan Iyer, and Yoav Artzi. 2018. Learning to map context-dependent sentences to executable formal queries. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 2238–2249, New Orleans, Louisiana. Association for Computational Linguistics.

Theodore R Sumers, Robert D Hawkins, Mark K Ho, and Thomas L Griffiths. 2021. Extending rational models of communication from beliefs to actions. In *Proceedings of the Annual Meeting of the Cognitive Science Society (CogSci)*.

Stefanie Tellex, Thomas Kollar, Steven Dickerson, Matthew R. Walter, Ashis Gopal Banerjee, Seth J. Teller, and Nicholas Roy. 2011. Understanding natural language commands for robotic navigation and mobile manipulation. In *Proceedings of the Twenty-Fifth AAAI Conference on Artificial Intelligence, AAAI 2011, San Francisco, California, USA, August 7-11, 2011*. AAAI Press.

Jesse Thomason, Shiqi Zhang, Raymond J. Mooney, and Peter Stone. 2015. Learning to interpret natural language commands through human-robot dialog. In *Proceedings of the Twenty-Fourth International Joint Conference on Artificial Intelligence, IJCAI 2015, Buenos Aires, Argentina, July 25-31, 2015*, pages 1923–1929. AAAI Press.

Sida I. Wang, Percy Liang, and Christopher D. Manning. 2016. Learning language games through interaction. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2368–2378, Berlin, Germany. Association for Computational Linguistics.

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, and Jamie Brew. 2019. Huggingface's transformers: State-of-the-art natural language processing. *CoRR*, abs/1910.03771.

Brian D. Ziebart, Andrew Maas, J. Andrew Bagnell, and Anind K. Dey. 2008. Maximum entropy inverse reinforcement learning. In *Proceedings of the Twenty-Third AAAI Conference on Artificial Intelligence*, pages 1433–1438. AAAI Press.

## A  Data Collection

We recruited Amazon Mechanical Turk workers from the US with an approval rate of $\geq 98\%$, $\geq 5{,}000$ Human Intelligence Tasks (HITs) completed, and completion of a custom qualification task where they played 15 minutes of the game and demonstrated active participation from manual review. Turk workers were given the following instructions for the task (shortened):
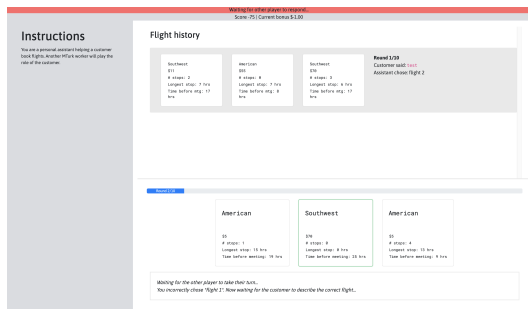
*Scenario: A customer working with a new personal assistant to book flights for their business meetings (you will play either the customer or the assistant, and another Turker will play the other role).*

*As the customer, you have specific preferences for what kind of flights you like and dislike. When you first start working with your assistant, you might need to tell them exactly what you want, but you hope that over time, they will figure out what you like and book your preferred flights without your help (imagine an assistant that knows you well enough to say: "Bob hates redeyes and doesn't like to be rushed on the way to the airport, so I'll go ahead and book this 2pm Jetblue to New York.")*

*As the assistant, you want to figure out what the customer likes and book flights that they want. Pay attention to what they say when they choose flights.*

*What is this task for? The goal of this task is to study how people naturally communicate their preferences to personal assistants, with the goal of building digital assistants (like Siri) that better understand and learn what people want. For the purposes of the task, we will give the customer "fake" flight preferences. If you are the customer, pretend that these are the kinds of flights you actually like / dislike.*

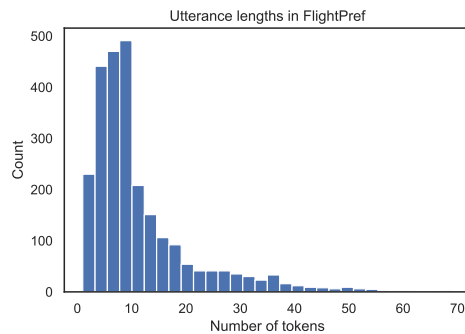The interface for the task (for an assistant) is shown below:



Workers were compensated $4 for 30 minutes (6 games of 6 rounds each), with a $0.20 bonus for every 25 points. For each round, points were accumulated based on the assistant's action:

- Assistant chooses correctly: +25 points

- Assistant chooses incorrectly: -100 points

- Assistant asks for help: -20 points

## B  Dataset Details

We plot the distribution of utterance lengths (number of tokens) in our dataset below:



## C  Base Model Details

### C.1  Training

We fine-tune all model parameters, including the parameters of the initially-pretrained BERT utterance encoders in the listener model and speaker model. We produce utterance representations from BERT using a linear projection of BERT's [CLS] embedding for the utterance. Models are trained separately on the training data using the loss functions below, with an AdamW learning rate of $2 \times 10^{-5}$ and a batch size of 64 for the listener model and a learning rate of $5 \times 10^{-5}$ and a batch size of 32 for the speaker model. We perform early stopping using the loss on held-out validation data.

**Listener loss.** We define the following loss function for the base listener model $L_{\text{base}}$ on each training instance:

$$\mathcal{L}_{L_{\text{base}}} = -\log p_{L_{\text{base}}}(\xi \mid u, \mathcal{M})$$

**Speaker loss.** It would be possible to learn the parameters of the base reward speaker $S_{\text{base}}$ directly on the training data, using a loss function similar to the base listener model. However, since utterances are often also action-descriptive, utterances cannot typically be predicted accurately from rewards alone. To account for this, we also define a separate base *action speaker* model, $p_{S_{\text{act}}}(u \mid \xi^*, \mathcal{M})$

that produces utterances conditioned on the optimal action $\xi^*$ in context $\mathcal{M}$:

$$p_{S_{\text{act}}}(u \mid \xi^*, \mathcal{M}) \propto \exp(\text{MLP}_{S_{\text{act}}}(\xi^*) \cdot \text{BERT}_S(u))$$

where $p_{S_{\text{act}}}$ is normalized over the same set of utterances as $p_{S_{\text{base}}}$. The base action speaker model is used only in training. It would be possible to also use this base action speaker in evaluation, in place of the pragmatic action speaker $p_{\text{action}}$ (Section 3.1); however, we found that pragmatic reasoning about a conventional instruction following model, as outlined in Section 3.1, performs better.

We train the two base speaker models jointly using a simple latent variable model, which makes the simplifying assumption that every utterance is either action- or reward-descriptive. To model this, we use a discrete latent variable $\lambda \in \{0, 1\}$:

$$p(u \mid \theta, \xi^*, \mathcal{M}, \lambda) = \lambda p_{S_{\text{base}}}(u \mid \theta) + (1 - \lambda)p_{S_{\text{act}}}(u \mid \xi^*, \mathcal{M})$$

The loss function for a single example is

$$\mathcal{L}_{S_{\text{base}}} = -\log \sum_{\lambda \in \{0,1\}} p(\lambda)p(u \mid \theta, \xi^*, \mathcal{M}, \lambda)$$

where $p(\lambda)$ gives the probability of an utterance being reward-descriptive or action-descriptive. We model this using $p(\lambda = 1) = \sigma(l)$, where $l$ is a parameter which is updated in training. Intuitively, the latent variable model learns soft clusters for action-descriptive and reward-descriptive utterances, with reward-descriptive utterances providing stronger supervision for the $S_{\text{base}}$ model.

**Speaker normalization.** In training, we compute the normalizers for the $S_{\text{base}}$ model using all utterances within the mini-batch, as well as up to 4 synthetically-constructed hard-negative examples defined by replacing attribute mentions within the true utterance (detected using exact word match) with alternative distractor attributes (e.g., replacing any occurrences of "jetblue" with one of "southwest", "american", or "delta", randomly sampled). We found that constructing hard-negatives in this way allowed us to train the base speaker models effectively despite using a fairly small dataset and small training batch sizes.

In evaluation, we compute normalizers for the $S_{\text{base}}$ model using a filtered set of all utterances from the training data that contain no more than 8 words and no digits. (We use a smaller set of normalizers in training time for efficiency reasons.)

## C.2 Hyperparameters

MLPs use fully-connected layers with ReLU non-linearities, and dropout applied to each hidden representation during training. We show hyperparameters for the models in Table 2. The BERT model is BERT-base, implemented in HuggingFace's Transformers library (Wolf et al., 2019).

| Listener Hyperparameters | |
|---|---|
| $\text{MLP}_{L_{\text{base}}}$ hidden layers | 2 |
| $\text{MLP}_{L_{\text{base}}}$ hidden size | 768 |
| $\text{MLP}_{L_{\text{base}}}$ output size | 768 |
| $\text{MLP}_{L_{\text{base}}}$ dropout | 0.1 |
| **Speaker Hyperparameters** | |
| $\text{MLP}_{S_{\text{base}}}$ hidden layers | 2 |
| $\text{MLP}_{S_{\text{base}}}$ hidden size | 512 |
| $\text{MLP}_{S_{\text{base}}}$ output size | 128 |
| $\text{MLP}_{S_{\text{base}}}$ dropout | 0.2 |

Table 2: Hyperparameters for the base speaker and listener models.

## D Analysis

### D.1 Analyzing the effect of pragmatically modeling the speaker

While we do not expect pragmatically modeling an action-only speaker will help in our domain since the action space is small (there is little ambiguity in what the referenced action is), we explore the effect of pragmatically modeling a belief-only speaker. We compare the belief-only model to two non-pragmatic alternatives that directly infer $p(\theta \mid u)$ without explicitly calculating the speaker's distribution over utterances $p(u \mid \theta)$: (1) inference: normalizing the logits of the $S_{\text{base}}$ model over rewards $\theta$ rather than utterances, and (2) training: a $S_{\text{base}}$ model trained to maximize $p(\theta \mid u)$ instead of $p(u \mid \theta)$. We show the results in Table 3: our model outperforms on held-out accuracy by 5-6% over non-pragmatic alternatives, suggesting that

| Method | Held-out accuracy (%) |
|---|---|
| Action + reward (Ours) | $59.1 \pm 0.96$ |
| $p(\theta \mid u)$ inference | $53.6 \pm 1.10$ |
| $p(\theta \mid u)$ training | $52.7 \pm 0.92$ |

Table 3: Average held-out accuracy (over 1000 rounds) of our model compared to ablated baselines that do not explicitly calculate a normalized utterance distribution $p(u \mid \theta)$, averaged over all validation rounds (each with varying numbers of observations). Standard error of the mean indicated and $p < .05$ for all observed differences using the paired bootstrap test.

```
Action-only better

+0.37: the cheaper the better
+0.27: like american airlines best
...
american
cheapest one please
...
-0.45: i love american and like southwest. i don't like
jetblue. i like low number of stops, but i like long stop
times.
-0.52: if i had a choice, i would never fly with delta and
american! get me jetblue or southwest if possible! if i
didn't have a choice, i really like having long stopovers,
so i can rest or sightsee. i also like having some time
before meetings so i'm not rushed.


Belief-only better
```

Figure 7: Utterances with the largest difference between action-only and belief-only held-out accuracy after updating independently on that round (compared to two random other utterances in the validation set), for examples where they both do better than chance (.33). The difference (action_acc - belief_acc) is shown to the left of each example.

modeling the speaker distribution is helpful for interpreting utterances more accurately.

## D.2 Examples ranked by action-only and belief-only accuracy difference

Figure 7 shows utterances with the largest difference between action-only and belief-only held-out accuracy after updating independently on that round (compared to two random other utterances in the validation set), for examples where they both do better than chance (.33). The belief-only model excels at reward-descriptive utterances, whereas action-only flight tends to outperform when there is less information in the utterance.