

Subgoal Discovery for Hierarchical Dialogue Policy Learning

Da Tang* Xiujun Li† Jianfeng Gao† Chong Wang‡ Lihong Li‡ Tony Jebara*

†Microsoft Research, Redmond, WA, USA

*Columbia University, NY, USA

‡Google Inc., Kirkland, WA, USA

{xiul, jfgao}@microsoft.com

{datang, jebara}@cs.columbia.edu

{chongw, lihong}@google.com

Abstract

Developing agents to engage in complex goal-oriented dialogues is challenging partly because the main learning signals are very sparse in long conversations. In this paper, we propose a divide-and-conquer approach that discovers and exploits the hidden structure of the task to enable efficient policy learning. First, given successful example dialogues, we propose the *Subgoal Discovery Network* (SDN) to divide a complex goal-oriented task into a set of simpler subgoals in an *unsupervised* fashion. We then use these subgoals to learn a multi-level policy by hierarchical reinforcement learning. We demonstrate our method by building a dialogue agent for the composite task of travel planning. Experiments with simulated and real users show that our approach performs competitively against a state-of-the-art method that requires human-defined subgoals. Moreover, we show that the learned subgoals are often human comprehensible.

1 Introduction

Consider we want to plan a trip to a distant city using a dialogue agent. The agent must make choices at each leg, e.g., whether to fly or to drive, whether to book a hotel. Each of these steps in turn involves making a sequence of decisions all the way down to *lower*-level actions. For example, to book a hotel involves identifying the location, specifying the check-in date and time, and negotiating the price etc.

The above process of the agent has a natural hierarchy: a top-level process selects which subgoal to complete, and a low-level process chooses primitive actions to accomplish the selected subgoal. Within the reinforcement learning (RL) paradigm, such a hierarchical decision making process can be formulated in the *options* framework (Sutton et al., 1999), where subgoals with

their own reward functions are used to learn policies for achieving these subgoals. These learned policies are then used as temporally extended actions, or options, for solving the entire task.

Based on the options framework, researchers have developed dialogue agents for complex tasks, such as travel planning, using hierarchical reinforcement learning (HRL) (Cuayáhuitl et al., 2010). Recently, Peng et al. (2017b) showed that the use of subgoals mitigates the reward sparsity and leads to more effective exploration for dialogue policy learning. However, these subgoals need to be human-defined which limits the applicability of the approach in practice because the domain knowledge required to properly define subgoals is often not available in many cases.

In this paper, we propose a simple yet effective *Subgoal Discovery Network* (SDN) that discovers useful subgoals automatically for an RL-based dialogue agent. The SDN takes as input a collection of successful conversations, and identifies “hub” states as subgoals. Intuitively, a hub state is a region in the agent’s state space that the agent tends to visit frequently on successful paths to a goal but not on unsuccessful paths. Given the discovered subgoals, HRL can be applied to learn a hierarchical dialogue policy which consists of (1) a top-level policy that selects among subgoals, and (2) a low-level policy that chooses primitive actions to achieve selected subgoals.

We present the first study of learning dialogue agents with automatically discovered subgoals. We demonstrate the effectiveness of our approach by building a composite task-completion dialogue agent for travel planning. Experiments with both simulated and real users show that an agent learned with discovered subgoals performs competitively against an agent learned using expert-defined subgoals, and significantly outperforms an agent learned without subgoals. We also find that

the subgoals discovered by SDN are often human comprehensible.

2 Background

A goal-oriented dialogue can be formulated as a Markov decision process, or MDP (Levin et al., 2000), in which the agent interacts with its environment over a sequence of discrete steps. At each step $t \in \{0, 1, \dots\}$, the agent observes the current state s_t of the conversation (Henderson, 2015; Mrkšić et al., 2017; Li et al., 2017), and chooses action a_t according to a policy π . Here, the action may be a natural-language sentence or a speech act, among others. Then, the agent receives a numerical reward r_t and switches to next state s_{t+1} . The process repeats until the dialogue terminates. The agent is to learn to choose *optimal* actions $\{a_t\}_{t=1,2,\dots}$ so as to maximize the total *discounted* reward $r_0 + \gamma r_1 + \gamma^2 r_2 + \dots$, where $\gamma \in [0, 1]$ is a discount factor. This learning paradigm is known as reinforcement learning, or RL (Sutton and Barto, 1998).

When facing a complex task, it is often more efficient to divide it into multiple simpler sub-tasks, solve them, and combine the partial solutions into a full solution for the original task. Such an approach may be formalized as hierarchical RL (HRL) in the options framework (Sutton et al., 1999). An option can be understood as a subgoal, which consists of an initiation condition (when the subgoal can be triggered), an option policy to solve the subgoal, and a termination condition (when the subgoal is considered finished).

When subgoals are given, there exist effective RL algorithms to learn a hierarchical policy. A major open challenge is the automatic discovery of subgoals from data, the main innovation of this work is covered in the next section.

3 Subgoal Discovery for HRL

Figure 1 shows the overall workflow of our proposed method of using automatic subgoal discovery for HRL. First a dialogue session is divided into several segments. Then at the end of those segments (subgoals), we equip an intrinsic or extrinsic reward for the HRL algorithm to learn a hierarchical dialogue policy. Note that only the last segment has an extrinsic reward. The details of the segmentation algorithm and how to use subgoals for HRL are presented in Section 3.1 and Section 3.3.



Figure 1: The workflow for HRL with subgoal discovery. In addition to the extrinsic reward at the end of the dialogue session, HRL also uses intrinsic rewards induced by the subgoals (or the ends of dialogue segments). Section 3.2 details the reward design for HRL.

3.1 Subgoal Discovery Network

Assume that we have collected a set of successful state trajectories of a task, as shown in Figure 2. We want to find subgoal states, such as the three red states s_4 , s_9 and s_{13} , which form the “hubs” of these trajectories. These hub states indicate the subgoals, and thus divide a state trajectory into several segments, each for an option¹.

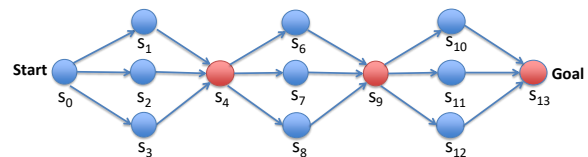


Figure 2: Illustration of “subgoals”. Assuming that there are three state trajectories $(s_0, s_1, s_4, s_6, s_9, s_{10}, s_{13})$, $(s_0, s_2, s_4, s_7, s_9, s_{11}, s_{13})$ and $(s_0, s_3, s_4, s_8, s_9, s_{12}, s_{13})$. Then red states s_4 , s_9 , s_{13} could be good candidates for “subgoals”.

Thus, discovering subgoals by identifying hubs in state trajectories is equivalent to segmenting state trajectories into options. In this work, we formulate subgoal discovery as a state trajectory segmentation problem, and address it using the Subgoal Discovery Network (SDN), inspired by the sequence segmentation model (Wang et al., 2017).

The SDN architecture. SDN repeats a two-stage process of generating a state trajectory segment, until a trajectory termination symbol is generated: first it uses an initial segment hidden state

¹There are many ways of creating a new option $\langle I, \pi, \beta \rangle$ for a discovered subgoal state. For example, when a subgoal state is identified at time step t , we add to I the set of states visited by the agent from time $t - n$ to t , where n is a pre-set parameter. I is therefore the union of all such states over all the state trajectories. The termination condition β is set to 1 when the subgoal is reached or when the agent is no longer in I , and to 0 otherwise. In the deep RL setting where states are represented by continuous vectors, β is a probability whose value is proportional to the vector distance e.g., between current state and subgoal state.

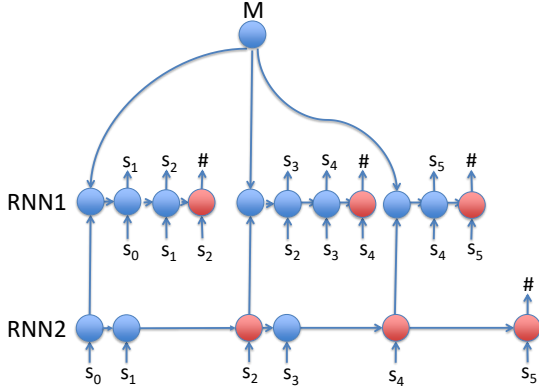


Figure 3: Illustration of SDN for state trajectory (s_0, \dots, s_5) with s_2 , s_4 and s_5 as subgoals. Symbol # is the termination. The top-level RNN (RNN1) models segments and the low-level RNN (RNN2) provides information about previous states from RNN1. The embedding matrix M maps the outputs of RNN2 to low dimensional representations so as to be consistent with the input dimensionality of RNN1. Note that state s_5 is associated with two termination symbols #; one is for the termination of the last segment and the other is for the termination of the entire trajectory.

to start a new segment, or a trajectory termination symbol to terminate the trajectory, given all previous states; if the trajectory is not terminated, then keep generating the next state in this trajectory segment given previous states until a segment termination symbol is generated. We illustrated this process in Figure 3.

We model the likelihood of each segment using an RNN, denoted as RNN1. During the training, at each time step, RNN1 predicts the next state with the current state as input, until it reaches the option termination symbol #. Since different options are under different conditions, it is not plausible to apply a fixed initial input to each segment. Therefore, we use another RNN (RNN2) to encode all previous states to provide relevant information and we transform these information to low dimensional representations as the initial inputs for the RNN1 instances. This is based on the *causality* assumption of the options framework (Sutton et al., 1999) — the agent should be able to determine the next option given all previous information, and this should not depend on information related to any later state. The low dimensional representations are obtained via a global subgoal embedding matrix $M \in \mathbb{R}^{d \times D}$, where d and D are the dimensionality of RNN1’s input layer and RNN2’s output layer, respectively. Mathematically, if the

output of RNN2 at time step t is o_t , then from time t the RNN1 instance has $M \cdot \text{softmax}(o_t)$ as its initial input². D is the number of subgoals we aim to learn. Ideally, the vector $\text{softmax}(o_t)$ in a well-trained SDN is close to an one-hot vector. Therefore, $M \cdot \text{softmax}(o_t)$ should be close to one column in M and we can view that M provides at most D different “embedding vectors” for RNN1 as inputs, indicating at most D different subgoals. Even in the case where $\text{softmax}(o_t)$ is not close to any one-hot vector, choosing a small D helps avoid overfitting.

Segmentation likelihood. Given the state trajectory (s_0, \dots, s_5) , assuming that s_2 , s_4 and s_5 are the discovered subgoal states, we model the conditional likelihood of a proposed segmentation $\sigma = ((s_0, s_1, s_2), (s_2, s_3, s_4), (s_4, s_5))$ as $p(\sigma|s_0) = p((s_0, s_1, s_2)|s_0) \cdot p((s_2, s_3, s_4)|s_{0:2}) \cdot p((s_4, s_5)|s_{0:4})$, where each probability term $p(\cdot|s_{0:i})$ is based on an RNN1 instance. And for the whole trajectory (s_0, \dots, s_5) , its likelihood is the sum over all possible segmentations.

Generally, for state trajectory $\mathbf{s} = (s_0, \dots, s_T)$, we model its likelihood as follows³:

$$L_S(\mathbf{s}) = \sum_{\sigma \subseteq \mathcal{S}(\mathbf{s}), \text{length}(\sigma) \leq S} \prod_{i=1}^{\text{length}(\sigma)} p(\sigma_i | \tau(\sigma_{1:i})), \quad (1)$$

where $\mathcal{S}(\mathbf{s})$ is the set of all possible segmentations for the trajectory \mathbf{s} , σ_i denotes the i^{th} segment in the segmentation σ , and τ is the concatenation operator. S is an upper limit on the maximal number of segments. This parameter is important for learning subgoals in our setting since we usually prefer a small number of subgoals. This is different from Wang et al. (2017), where a maximum segment length is enforced.

We use maximum likelihood estimation with Eq. (1) for training. However, the number of possible segmentations is exponential in $\mathcal{S}(\mathbf{s})$ and the naive enumeration is intractable. Here, dynamic programming is employed to compute the likelihood in Eq. (1) efficiently: for a trajectory $\mathbf{s} = (s_0, \dots, s_T)$, if we denote the sub-trajectory (s_i, \dots, s_t) of \mathbf{s} as $\mathbf{s}_{i:t}$, then its likelihood follows

² $\text{softmax}(o_t)_i = \exp(o_{t,i}) / \sum_{i'=1}^D \exp(o_{t,i'}) \in \mathbb{R}^D$ for $o_t = (o_{t,1}, \dots, o_{t,D})$.

³For notation convenience, we include s_0 into the observational sequence, though s_0 is always conditioned upon.

the below recursion:

$$L_m(\mathbf{s}_{0:t}) = \begin{cases} \sum_{i=0}^{t-1} L_{m-1}(\mathbf{s}_{0:i})p(\mathbf{s}_{i:t}|\mathbf{s}_{0:i}), & m > 0, \\ I[t = 0], & m = 0. \end{cases}$$

Here, $L_m(\mathbf{s}_{0:t})$ denotes the likelihood of sub-trajectory $\mathbf{s}_{0:t}$ with no more than m segments and $I[\cdot]$ is an indicator function. $p(\mathbf{s}_{i:t}|\mathbf{s}_{0:i})$ is the likelihood segment $\mathbf{s}_{i:t}$ given the previous history, where RNN1 models the segment and RNN2 models the history as shown in Figure 3. With this recursion, we can compute the likelihood $L_S(\mathbf{s})$ for the trajectory $\mathbf{s} = (s_0, \dots, s_T)$ in $O(ST^2)$ time.

Learning algorithm. We denote θ^s as the model parameter including the parameters of the embedding matrix M , RNN1 and RNN2. We then parameterize the segment likelihood function as $p(\mathbf{s}_{i:t}|\mathbf{s}_{0:i}) = p(\mathbf{s}_{i:t}|\mathbf{s}_{0:i}; \theta^s)$, and the trajectory likelihood function as $L_m(\mathbf{s}_{0:t}) = L_m(\mathbf{s}_{0:t}; \theta^s)$.

Given a set of N state trajectories $(\mathbf{s}^{(1)}, \dots, \mathbf{s}^{(N)})$, we optimize θ^s by minimizing the negative mean log-likelihood with L_2 regularization term $\frac{1}{2}\lambda\|\theta^s\|^2$ where $\lambda > 0$, using stochastic gradient descent:

$$\mathcal{L}_S(\theta^s, \lambda) = -\frac{1}{N} \sum_{i=1}^N \log L_S(\mathbf{s}^{(i)}, \theta^s) + \frac{1}{2}\lambda\|\theta^s\|^2. \quad (2)$$

Algorithm 1 outlines the training procedure for SDN using stochastic gradient descent.

Algorithm 1 Learning SDN

Input: A set of state trajectories $(\mathbf{s}_1, \dots, \mathbf{s}_N)$, the number of segments limit S , initial learning rate $\eta > 0$.

- 1: Initialize the SDN parameter θ^s .
 - 2: **while** not converged **do**
 - 3: Compute the gradient $\nabla_{\theta^s} \mathcal{L}_S(\theta^s, \lambda)$ of the loss $\mathcal{L}_S(\theta^s, \lambda)$ as in Eq. (2).
 - 4: Update $\theta^s \leftarrow \theta^s - \eta \nabla_{\theta^s} \mathcal{L}_S(\theta^s, \lambda)$.
 - 5: Update the learning rate η .
 - 6: **end while**
-

3.2 Hierarchical Dialogue Policy Learning

Before describing how we use a trained SDN model for HRL, we first present a short review of HRL for a task-oriented dialogue system. Following the *options* framework (Sutton et al., 1999), assume that we have a state set \mathcal{S} , an option set \mathcal{G} and a finite primitive action set \mathcal{A} .

The HRL approach we take learns two Q-functions (Peng et al., 2017b), parameterized by θ_e and θ_i , respectively:

- The top-level $Q^*(s, g; \theta_e)$ measures the maximum total discounted *extrinsic* reward received by choosing subgoal g in state s and then following an optimal policy. These extrinsic rewards are the objective to be maximized by the entire dialogue policy.
- The low-level $Q^*(s, a, g; \theta_i)$ measures the maximum total discounted intrinsic reward received to achieve a *given* subgoal g , by choosing action a in state s and then following an optimal option policy. These intrinsic rewards are used to learn an option policy to achieve a given subgoal.

Suppose we have a dialogue session of T turns: $\tau = (s_0, a_0, r_0, \dots, s_T)$, which is segmented into a sequence of subgoals $g_0, g_1, \dots \in \mathcal{G}$. Consider one of these subgoals g which starts and ends in steps t_0 and t_1 , respectively.

The top-level Q-function is learned using Q-learning, by treating subgoals as temporally extended actions:

$$\theta_e \leftarrow \theta_e + \alpha \cdot (q - Q(s_t, g; \theta_e)) \cdot \nabla_{\theta_e} Q(s_t, g; \theta_e),$$

where

$$q = \sum_{t=t_0}^{t_1-1} \gamma^{t-t_0} r_t^e + \gamma^{t_1-t_0} \max_{g' \in \mathcal{G}} Q(s_{t_1}, g'; \theta_e),$$

and α is the step-size parameter, $\gamma \in [0, 1]$ is a discount factor. In the above expression of q , the first term refers to the total discounted reward during fulfillment of subgoal g , and the second to the maximum total discounted after g is fulfilled.

The low-level Q-function is learned in a similar way, and follows the standard Q-learning update, except that intrinsic rewards for subgoal g are used. Specifically, for $t = t_0, t_0 + 1, \dots, t_1 - 1$:

$$\theta_i \leftarrow \theta_i + \alpha \cdot (q_t - Q(s_t, a_t, g; \theta_e)) \cdot \nabla_{\theta_i} Q(s_t, a_t, g; \theta_i),$$

where

$$q_t = r_t^i + \gamma \max_{a' \in \mathcal{A}} Q(s_{t+1}, a', g; \theta_i).$$

Here, the intrinsic reward r_t^i is provided by the internal critic of dialogue manager. More details are in Appendix A.

In hierarchical policy learning, the combination of the extrinsic and intrinsic rewards is expected to help the agent to successfully accomplish a composite task as fast as possible while trying to avoid unnecessary subtask switches. Hence, we define the extrinsic and intrinsic rewards as follows:

Extrinsic Reward. Let L be the maximum number of turns of a dialogue, and K the number of subgoals. At the end of a dialogue, the agent receives a positive extrinsic reward of $2L$ for a success dialogue, or $-L$ for a failure dialogue; for each turn, the agent receives an extrinsic reward of -1 to encourage shorter dialogues.

Intrinsic Reward. When a subgoal terminates, the agent receives a positive intrinsic reward of $2L/K$ if a subgoal is completed successfully, or a negative intrinsic reward of -1 otherwise; for each turn, the agent receives an intrinsic reward -1 to encourage shorter dialogues.

3.3 Hierarchical Policy Learning with SDN

We use a trained SDN in HRL as follows. The agent starts from the initial state s_0 , keeps sampling the output from the distribution related to the top-level RNN (RNN1) until a termination symbol $\#$ is generated, which indicates the agent reaches a subgoal. In this process, intrinsic rewards are generated as specified in the previous subsection. After $\#$ is generated, the agent selects a new option, and repeats this process.

This type of naive sampling may allow the option to terminate at some places with a low probability. To stabilize the HRL training, we introduce a threshold $p \in (0, 1)$, which directs the agent to terminate an option if and only if the probability of outputting $\#$ is at least p . We found this modification leads to better behavior of the HRL agent than the naive sampling method, since it normally has a smaller variance.

In the HRL training, the agent only uses the probability of outputting $\#$ to decide subgoal termination. Algorithm 2 outlines the full procedure of one episode for hierarchical dialogue policies with a trained SDN in the composite task-completion dialogue system.

4 Experiments and Results

We evaluate the proposed model on a travel planning scenario for composite task-oriented dialogues (Peng et al., 2017b). Over the exchange of a conversation, the agent gathers information about the user’s intent before booking a trip. The environment then assesses a binary outcome (success or failure) at the end of the conversation, based on (1) whether a trip is booked, and (2) whether the trip satisfies the user’s constraints.

Algorithm 2 HRL episode with a trained SDN

Input: A trained SDN \mathcal{M} , initial state s_0 of an episode, threshold p , the HRL agent \mathcal{A} .

- 1: Initialize an RNN2 instance R_2 with parameters from \mathcal{M} and s_0 as the initial input.
- 2: Initialize an RNN1 instance R_1 with parameters from \mathcal{M} and $M \cdot \text{softmax}(o_0^{\text{RNN2}})$ as the initial input, where M is the embedding matrix (from \mathcal{M}) and o_0^{RNN2} is the initial output of R_2 .
- 3: Current state $s \leftarrow s_0$.
- 4: Select an option o using the agent \mathcal{A} .
- 5: **while** Not reached the final goal **do**
- 6: Select an action a according to s and o using the agent \mathcal{A} . Get the reward r and the next state s' from the environment.
- 7: Place s' to R_2 , denote o_t^{RNN2} as R_2 ’s latest output and take $M \cdot \text{softmax}(o_t^{\text{RNN2}})$ as the R_1 ’s new input. Let $p_{s'}$ be the probability of outputting the termination symbol $\#$.
- 8: **if** $p_{s'} \geq p$ **then**
- 9: Select a new option o using the agent \mathcal{A} .
- 10: Re-initialize R_1 using the latest output from R_2 and the embedding matrix M .
- 11: **end if**
- 12: **end while**

Dataset. The raw dataset in our experiments is from a publicly available multi-domain dialogue corpus (El Asri et al., 2017). Following Peng et al. (2017b), a few changes were made to introduce dependencies among subtasks. For example, the hotel check-in date should be the same with the departure flight arrival date. The data was mainly used to create simulated users, and to build the knowledge bases for the subtasks of booking flights and reserving hotels.

User Simulator. In order to learn good policies, RL algorithms typically need an environment to interact with. In the dialogue research community, it is common to use simulated users for this purpose (Schatzmann et al., 2007; Li et al., 2017; Liu and Lane, 2017). In this work, we adapted a publicly available user simulator (Li et al., 2016) to the composite task-completion dialogue setting with the dataset described above. During training, the simulator provides the agent with an (extrinsic) reward signal at the end of the dialogue. A dialogue is considered to be successful only when a travel plan is booked successfully, and the information provided by the agent satisfies user’s constraints.

Baseline Agents. We benchmarked the proposed agent (referred to as the *m-HRL Agent*) against three baseline agents:

- A *Rule Agent* uses a sophisticated, hand-crafted dialogue policy, which requests and informs a hand-picked subset of necessary slots, and then

confirms with the user about the reserved trip before booking the flight and hotel.

- A *flat RL Agent* is trained with a standard deep reinforcement learning method, DQN (Mnih et al., 2015), which learns a flat dialogue policy using extrinsic rewards only.
- A *h-HRL Agent* is trained with hierarchical deep reinforcement learning (HDQN), which learns a hierarchical dialogue policy based on human-defined subgoals (Peng et al., 2017b).

Collecting State Trajectories. Recall that our subgoal discovery approach takes as input a set of state trajectories which lead to successful outcomes. In practice, one can collect a large set of successful state trajectories, either by asking human experts to demonstrate (e.g., in a call center), or by rolling out a reasonably good policy (e.g., a policy designed by human experts). In this paper, we obtain dialogue state trajectories from a rule-based agent which is handcrafted by a domain expert, the performance of this rule-based agent can achieve success rate of 32.2% as shown in Figure 4 and Table 1. We only collect the successful dialogue sessions from the roll-outs of the rule-based agent, and try to learn the subgoals from these dialogue state trajectories.

Experiment Settings. To train SDN, we use RMSProp (Tieleman and Hinton, 2012) to optimize the model parameters. For both RNN1 and RNN2, we use LSTM (Hochreiter and Schmidhuber, 1997) as hidden units and set the hidden size to 50. We set embedding matrix M with $D = 4$ columns. As we discussed in Section 3.1, D captures the maximum number of subgoals that the model is expected to learn. Again, to avoid SDN from learning many unnecessary subgoals, we only allow segmentation with at most $S = 4$ segments during subgoal training. The values for D and S are usually set to be a little bit larger than the expected number of subgoals (e.g., 2 or 3 for this task) since we expect a great proportion of the subgoals that SDN learns are useful, but not necessary for all of them. As long as SDN discovers useful subgoals that guide the agent to learn policies faster, it is beneficial for HRL training, even if some non-perfect subgoals are found. During the HRL training, we use the learned SDN to propose subgoal-completion queries. In our experiment, we set the maximum turn $L = 60$.

We collected $N = 1634$ successful, but imper-

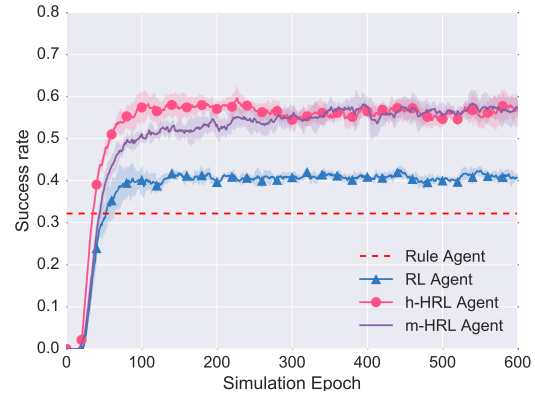


Figure 4: Learning curves of agents under simulation.

Agent	Success Rate	Turns	Reward
Rule	.3220	46.23	-24.02
RL	.4440	45.50	-1.834
h-HRL	.6485	44.23	35.32
m-HRL	.6455	44.85	34.77

Table 1: Performance of agents with simulated user.

fect, dialogue episodes from the rule-based agent in Table 1 and randomly choose 80% of these dialogue state trajectories for training SDN. The remaining 20% were used as a validation set.

As illustrated in Section 3.3, SDN starts a new RNN1 instance and issues a subgoal-completion query when the probability of outputting the termination symbol # is above a certain threshold p (as in Algorithm 2). In our experiment, p is set to be 0.2, which was manually picked according to the termination probability during SDN training.

In dialogue policy learning, for the baseline RL agent, we set the size of the hidden layer to 80. For the HRL agents, both top-level and low-level dialogue policies have a hidden layer size of 80. RMSprop was applied to optimize the parameters. We set the batch size to be 16. During training, we used ϵ -greedy strategy for exploration with annealing and set $\gamma = 0.95$. For each simulation epoch, we simulated 100 dialogues and stored these state transition tuples in the experience replay buffers. At the end of each simulation epoch, the model was updated with all the transition tuples in the buffers in a batch manner.

4.1 Simulated User Evaluation

In the composite task-completion dialogue scenario, we compared the proposed *m-HRL* agent

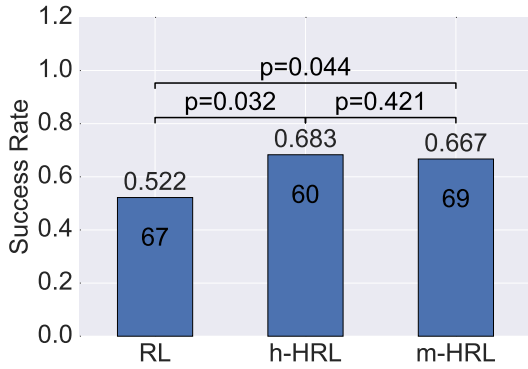


Figure 5: Performance of three agents tested with real users: success rate, number of dialogues and p-value are indicated on each bar (difference in mean is significant with $p < 0.05$).

with three baseline agents in terms of three metrics: success rate⁴, average rewards and average turns per dialogue session.

Figure 4 shows the learning curves of all four agents trained against the simulated user. Each learning curve was averaged over 5 runs. Table 1 shows the test performance where each number was averaged over 5 runs and each run generated 2000 simulated dialogues. We find that the HRL agents generated higher success rates and needed fewer conversation turns to achieve the users’ goals than the rule-based agent and the flat RL agent. The performance of the m-HRL agent is tied with that of the h-HRL agent, even though the latter requires high-quality subgoals designed by human experts.

4.2 Human Evaluation

We further evaluated the agents that were trained on simulated users against real users, who were recruited from the authors’ organization. We conducted a study using the one RL agent and two HRL agents {RL, h-HRL, m-HRL}, and compared two pairs: {RL, m-HRL} and {h-HRL, m-HRL}. In each dialogue session, one agent was randomly selected from the pool to interact with a user. The user was *not* aware of which agent was selected to avoid systematic bias. The user was presented with a goal sampled from a user-goal corpus, then was instructed to converse with the agent to complete the given task. At the end of each dialogue session, the user was asked to give a rating on a scale from 1 to 5 based on the natural-

⁴Success rate is the fraction of dialogues which accomplished the task successfully within the maximum turns.

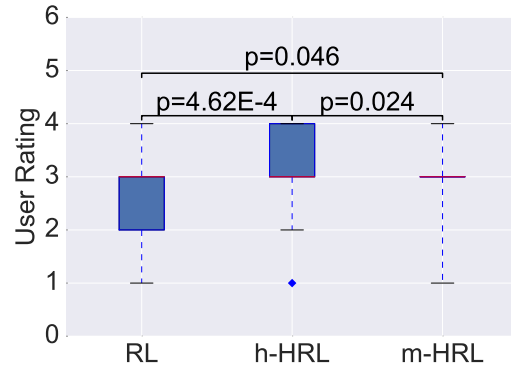


Figure 6: Distribution of user ratings for three agents in human evaluation

ness and coherence of the dialogue; here, 1 is the worst rating and 5 the best. In total, we collected 196 dialogue sessions from 10 human users.

Figure 5 summarizes the performances of these agents against real users in terms of success rate. Figure 6 shows the distribution of user ratings for each agent. For these two metrics, both HRL agents were significantly better than the flat RL agent. Another interesting observation is that the m-HRL agent performs similarly to the h-HRL agent in terms of success rate in the real user study as shown on Figure 5. Meanwhile in Figure 6, the h-HRL agent is significantly better than m-HRL agent in terms of real user ratings. This may be caused by the probabilistic termination of subgoals: we used a threshold strategy to decide whether to terminate a subgoal. This could introduce variance so the agent might not behave reasonably compared with human-defined subgoals which terminate deterministically.

4.3 Subgoal Visualization

Table 2 shows the subgoals discovered by SDN in a sample dialogue by a rule-based agent interacting with the simulated user. The rule-based agent is equipped with a human-defined subtask structure, which always solves subtask *flight* (turns 1–15) before *hotel* (turns 16–23), as shown in the first column. At turn 10, the user starts to talk about *hotel* while the rule-based agent is still working on the pre-defined, unfinished *flight* subtask until subtask *flight* is finished at turn 15. At turn 16, the user switches to *hotel*, and so does the rule-based agent until the end of the dialogue. For this rule-based agent, the human-defined subgoal (*flight*) terminates at turn 15. Meanwhile, our SDN model detected two subgoals (except for the final goal):

h-Task	ID	Text
Flight	1	<i>agt</i> : Hi, may I know your departure city?
	2	<i>usr</i> : I depart from Fort Lauderdale.
	3	<i>agt</i> : And what is your destination city?
	4	<i>usr</i> : I need economy seat to Pittsburgh.
	5	<i>agt</i> : What is your departure date?
	6	<i>usr</i> : 2016-09-12, what time is available for the flight?
	7	<i>agt</i> : There is one flight at 21:00:00.
	8	<i>usr</i> : What is the price for that flight?
	9	<i>agt</i> : \$423. Subgoal 1: A detected subgoal terminates!
	10	<i>usr</i> : I need to reserve a hotel for five people.
	11	<i>agt</i> : What is your return date for the trip?
	12	<i>usr</i> : 2016-09-18, I need to book 5 tickets for the flight.
	13	<i>agt</i> : What time do you prefer for return flight?
	14	<i>usr</i> : What time is available that day?
	15	<i>agt</i> : There is one flight departing at 11:00:00AM. Subgoal 2: A detected subgoal terminates!
Hotel	16	<i>usr</i> : What hotel is available?
	17	<i>agt</i> : There is one hotel called Winter Hill Inn available.
	18	<i>usr</i> : What is the price for hotel?
	19	<i>agt</i> : \$410.
	20	<i>usr</i> : Great, could you book the trip for me?
	21	<i>agt</i> : Okay. Trip tickets are issued.
	22	<i>usr</i> : Thanks very much!
	23	<i>agt</i> : You are welcome!

Table 2: Discovered subgoals (except for the final goal) in a sample dialogue by a rule-based agent interacting with user simulator. The left column (h-Task) shows the human-defined subtasks for the rule-based agent. SDN detects two subgoals that terminate at turn 9 and 15 respectively. (h-Task: human-defined subtask, ID: turn ID, *agt*: Agent, *usr*: User)

one terminating at turn 9 (Subgoal 1), and another terminating at turn 15 (Subgoal 2). Subgoal 2 is consistent with the human-defined subgoal. Subgoal 1 is also reasonable since the user tries to switch to *hotel* at turn 10. In Appendix B, Table 3 shows a sample dialogue session by m-HRL agent interacting with a real user.

5 Related Work

Task-completion dialogue systems have attracted numerous research efforts, and there is growing interest in leveraging reinforcement learning for policy learning. One line of research is on single-domain task-completion dialogues with flat deep reinforcement learning algorithms such as DQN (Zhao and Eskenazi, 2016; Li et al., 2017; Peng et al., 2018), actor-critic (Peng et al., 2017a; Liu and Lane, 2017) and policy gradients (Williams et al., 2017; Liu et al., 2017). Another line of research addresses multi-domain dialogues where each domain is handled by a separate agent (Gašić et al., 2015; Gašić et al., 2015; Cuayáhuitl et al., 2016). Recently, Peng et al. (2017b) presented a composite task-completion dialogue system. Unlike multi-domain dialogue systems, composite tasks introduce inter-subtask constraints. As a result, the completion of a set

of individual subtasks does *not* guarantee the solution of the entire task.

Cuayáhuitl et al. (2010) applied HRL to dialogue policy learning, although they focus on problems with a small state space. Later, Budzianowski et al. (2017) used HRL in multi-domain dialogue systems. Peng et al. (2017b) first presented an HRL agent with a global state tracker to learn the dialogue policy in the composite task-completion dialogue systems. All these works are built based on subgoals that were pre-defined with human domain knowledge for the specific tasks. The only job of the policy learner is to learn a hierarchical dialogue policy, which leaves the subgoal discovery problem unsolved. In addition to the applications in dialogue systems, subgoal is also widely studied in the linguistics research community (Allwood, 2000; Linell, 2009).

In the literature, researchers have proposed algorithms to automatically discovery subgoals for hierarchical RL. One large body of work is based on analyzing the spatial structure of the state transition graphs, by identifying bottleneck states or clusters, among others (Stolle and Precup, 2002; McGovern and Barto, 2001; Mannor et al., 2004; Şimşek et al., 2005; Entezari et al., 2011; Bacon, 2013). Another family of algorithms identifies commonalities of policies and extracts these partial policies as useful skills (Thrun and Schwartz, 1994; Pickett and Barto, 2002; Brunskill and Li, 2014). While similar in spirit to ours, these methods do not easily scale to continuous problems as in dialogue systems. More recently, researchers have proposed deep learning models to discover subgoals in continuous-state MDPs (Bacon et al., 2017; Machado et al., 2017; Vezhnevets et al., 2017). It would be interesting to see how effective they are for dialogue management.

Segmental structures are common in human languages. In the NLP community, some related research on segmentation includes word segmentation (Gao et al., 2005; Zhang et al., 2016) to divide the words into meaningful units. Alternatively, topic detection and tracking (Allan et al., 1998; Sun et al., 2007) segment a stream of data and identify stories or events in news or social text. In this work, we formulate subgoal discovery as a trajectory segmentation problem. Section 3.1 presents our approach to subgoal discovery which is inspired by a probabilistic sequence segmentation model (Wang et al., 2017).

6 Discussion and Conclusion

We have proposed the Subgoal Discovery Network to learn subgoals automatically in an unsupervised fashion without human domain knowledge. Based on the discovered subgoals, we learn the dialogue policy for complex task-completion dialogue agents using HRL. Our experiments with both simulated and real users on a composite task of travel planning, show that an agent trained with automatically discovered subgoals performs competitively against an agent with human-defined subgoals, and significantly outperforms an agent without subgoals. Through visualization, we find that SDN discovers reasonable, comprehensible subgoals given only a small amount of suboptimal but successful dialogue state trajectories.

These promising results suggest several directions for future research. First, we want to integrate subgoal discovery into dialogue policy learning rather than treat them as two separate processes. Second, we would like to extend SDN to identify multi-level hierarchical structures among subgoals so that we can handle more complex tasks than those studied in this paper. Third, we would like to generalize SDN to a wide range of complex goal-oriented tasks beyond dialogue, such as the particularly challenging Atari game of Montezuma’s Revenge (Kulkarni et al., 2016).

Acknowledgments

We would like to thank the anonymous reviewers, members of the xlab at the University of Washington, and Chris Brockett, Michel Galley for their insightful comments on the work. Most of this work was done while DT, CW & LL were with Microsoft.

References

James Allan, Jaime G Carbonell, George Doddington, Jonathan Yamron, and Yiming Yang. 1998. Topic detection and tracking pilot study final report.

Jens Allwood. 2000. An activity based approach to pragmatics. *Abduction, belief and context in dialogue: Studies in computational pragmatics*, pages 47–80.

Pierre-Luc Bacon. 2013. On the bottleneck concept for options discovery: Theoretical underpinnings and extension in continuous state spaces. Master’s thesis, McGill University.

Pierre-Luc Bacon, Jean Harb, and Doina Precup. 2017. The option-critic architecture. In *Proceedings of the 31st AAAI Conference on Artificial Intelligence (AAAI)*, pages 1726–1734.

Emma Brunskill and Lihong Li. 2014. PAC-inspired option discovery in lifelong reinforcement learning. In *Proceedings of the 31st International Conference on Machine Learning (ICML)*, pages 316–324.

Pawel Budzianowski, Stefan Ultes, Pei-Hao Su, Nikola Mrksic, Tsung-Hsien Wen, Inigo Casanueva, Lina Rojas-Barahona, and Milica Gasic. 2017. Sub-domain modelling for dialogue management with hierarchical reinforcement learning. *arXiv preprint arXiv:1706.06210*.

Heriberto Cuayáhuitl, Steve Renals, Oliver Lemon, and Hiroshi Shimodaira. 2010. Evaluation of a hierarchical reinforcement learning spoken dialogue system. *Computer Speech and Language*, 24(2):395–429.

Heriberto Cuayáhuitl, Seunghak Yu, Ashley Williamson, and Jacob Carse. 2016. Deep reinforcement learning for multi-domain dialogue systems. *arXiv preprint arXiv:1611.08675*.

Layla El Asri, Hannes Schulz, Shikhar Sharma, Jeremie Zumer, Justin Harris, Emery Fine, Rahul Mehrotra, and Kaheer Suleman. 2017. Frames: A corpus for adding memory to goal-oriented dialogue systems. *arXiv:1704.00057*. URL: <https://datasets.maluuba.com/Frames>.

Negin Entezari, Mohammad Ebrahim Shiri, and Parham Moradi. 2011. Subgoal discovery in reinforcement learning using local graph clustering.

Jianfeng Gao, Mu Li, Andi Wu, and Chang-Ning Huang. 2005. Chinese word segmentation and named entity recognition: A pragmatic approach. *Computational Linguistics*, 31(4):531–574.

Milica Gašić, Dongho Kim, Pirros Tsiakoulis, and Steve Young. 2015. Distributed dialogue policies for multi-domain statistical dialogue management. In *ICASSP 2015*, pages 5371–5375.

Milica Gašić, Nikola Mrkšić, Pei hao Su, David Vandyke, Tsung-Hsien Wen, and Steve J. Young. 2015. Policy committee for adaptation in multi-domain spoken dialogue systems. In *ASRU*, pages 806–812. IEEE.

Matthew Henderson. 2015. Machine learning for dialog state tracking: A review. In *Proceedings of the 1st International Workshop on Machine Learning in Spoken Language Processing*.

Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation*, 9(8):1735–1780.

- Tejas D Kulkarni, Karthik Narasimhan, Ardavan Saedi, and Josh Tenenbaum. 2016. Hierarchical deep reinforcement learning: Integrating temporal abstraction and intrinsic motivation. In *Advances in Neural Information Processing Systems*, pages 3675–3683.
- Esther Levin, Roberto Pieraccini, and Wieland Eckert. 2000. A stochastic model of human-machine interaction for learning dialog strategies. *IEEE Transactions on Speech and Audio Processing*, 8(1):11–23.
- Xiujun Li, Zachary C Lipton, Bhuwan Dhingra, Lihong Li, Jianfeng Gao, and Yun-Nung Chen. 2016. A user simulator for task-completion dialogues. *arXiv preprint arXiv:1612.05688*.
- Xuijun Li, Yun-Nung Chen, Lihong Li, and Jianfeng Gao. 2017. End-to-end task-completion neural dialogue systems. *arXiv preprint arXiv:1703.01008*.
- Per Linell. 2009. *Rethinking Language, Mind, and World Dialogically*. Information Age Publishing.
- Bing Liu and Ian Lane. 2017. Iterative policy learning in end-to-end trainable task-oriented neural dialog models. *arXiv preprint arXiv:1709.06136*.
- Bing Liu, Gokhan Tur, Dilek Hakkani-Tur, Pararth Shah, and Larry Heck. 2017. End-to-end optimization of task-oriented dialogue model with deep reinforcement learning. *arXiv preprint arXiv:1711.10712*.
- Marlos C. Machado, Marc G. Bellemare, and Michael H. Bowling. 2017. A Laplacian framework for option discovery in reinforcement learning. In *Proceedings of the 34th International Conference on Machine Learning (ICML)*, pages 2295–2304.
- Shie Mannor, Ishai Menache, Amit Hoze, and Uri Klein. 2004. Dynamic abstraction in reinforcement learning via clustering. In *Proceedings of the 21st International Conference on Machine Learning (ICML)*, pages 560–567. ACM.
- Amy McGovern and Andrew G Barto. 2001. Automatic discovery of subgoals in reinforcement learning using diverse density. In *Proceedings of the 18th International Conference on Machine Learning*, pages 361–368.
- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. 2015. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529–533.
- Nikola Mrkšić, Diarmuid Ó Séaghdha, Tsung-Hsien Wen, Blaise Thomson, and Steve J. Young. 2017. Neural belief tracker: Data-driven dialogue state tracking. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics*, pages 1777–1788.
- Baolin Peng, Xiujun Li, Jianfeng Gao, Jingjing Liu, Yun-Nung Chen, and Kam-Fai Wong. 2017a. Adversarial advantage actor-critic model for task-completion dialogue policy learning. *arXiv preprint arXiv:1710.11277*.
- Baolin Peng, Xiujun Li, Jianfeng Gao, Jingjing Liu, and Kam-Fai Wong. 2018. Integrating planning for task-completion dialogue policy learning. *arXiv preprint arXiv:1801.06176*.
- Baolin Peng, Xiujun Li, Lihong Li, Jianfeng Gao, Asli Celikyilmaz, Sungjin Lee, and Kam-Fai Wong. 2017b. Composite task-completion dialogue policy learning via hierarchical deep reinforcement learning. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 2221–2230.
- Marc Pickett and Andrew G. Barto. 2002. Policy-Blocks: An algorithm for creating useful macro-actions in reinforcement learning. In *Proceedings of the 19th International Conference on Machine Learning (ICML)*, pages 506–513.
- Jost Schatzmann, Blaise Thomson, Karl Weilhammer, Hui Ye, and Steve Young. 2007. Agenda-based user simulation for bootstrapping a pomdp dialogue system. In *NAACL 2007; Companion Volume, Short Papers*, pages 149–152. Association for Computational Linguistics.
- Özgür Şimşek, Alicia P. Wolfe, and Andrew G. Barto. 2005. Identifying useful subgoals in reinforcement learning by local graph partitioning. In *Proceedings of the 22nd International Conference on Machine Learning (ICML)*, pages 816–823.
- Martin Stolle and Doina Precup. 2002. Learning options in reinforcement learning. In *International Symposium on Abstraction, Reformulation, and Approximation*, pages 212–223. Springer.
- Bingjun Sun, Prasenjit Mitra, C Lee Giles, John Yen, and Hongyuan Zha. 2007. Topic segmentation with shared topic detection and alignment of multiple documents. In *SIGIR 2007*, pages 199–206. ACM.
- Richard S. Sutton and Andrew G. Barto. 1998. *Reinforcement Learning: An Introduction*. MIT Press, Cambridge, MA.
- Richard S Sutton, Doina Precup, and Satinder Singh. 1999. Between MDPs and semi-MDPs: A framework for temporal abstraction in reinforcement learning. *Artificial intelligence*, 112(1-2):181–211.
- Sebastian Thrun and Anton Schwartz. 1994. Finding structure in reinforcement learning. In *Advances in Neural Information Processing Systems 7 (NIPS)*, pages 385–392.
- Tijmen Tieleman and Geoffrey Hinton. 2012. Lecture 6.5-rmsprop: Divide the gradient by a running average of its recent magnitude. *COURSERA: Neural networks for machine learning*, 4(2):26–31.

Alexander Sasha Vezhnevets, Simon Osindero, Tom Schaul, Nicolas Heess, Max Jaderberg, David Silver, and Koray Kavukcuoglu. 2017. FeUdal Networks for hierarchical reinforcement learning. In *Proceedings of the 34th International Conference on Machine Learning (ICML)*, pages 3540–3549.

Chong Wang, Yining Wang, Po-Sen Huang, Abdelrahman Mohamed, Dengyong Zhou, and Li Deng. 2017. Sequence modeling via segmentations. *arXiv preprint arXiv:1702.07463*.

Jason D Williams, Kavosh Asadi, and Geoffrey Zweig. 2017. Hybrid code networks: Practical and efficient end-to-end dialog control with supervised and reinforcement learning. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (ACL)*.

Meishan Zhang, Yue Zhang, and Guohong Fu. 2016. Transition-based neural word segmentation. In *ACL (1)*.

Tiancheng Zhao and Maxine Eskenazi. 2016. Towards end-to-end learning for dialog state tracking and management using deep reinforcement learning. *arXiv preprint arXiv:1606.02560*.

A Hierarchical Dialogue Policy Learning

This section provides more algorithmic details for Section 3.2. Again, assume a conversation of length T :

$$\tau = (s_0, a_0, r_0, \dots, s_{T-1}, a_{T-1}, r_{T-1}, s_T).$$

Suppose an HRL agent segments the trajectory into a sequence of subgoals as $g_0, g_1, \dots \in \mathcal{G}$, and the corresponding subgoal termination time steps as $t_{g_0}, t_{g_1}, \dots \in \mathbb{N}^*$. Furthermore, denote the intrinsic reward at time step t by r_t^i . The top-level and low-level Q-functions satisfy the following Bellman equations:

$$Q^*(s, g) = \mathbb{E} \left[\sum_{i=t_{g_j}}^{t_{g_{j+1}}-1} \gamma^{i-t_{g_j}} r_{t+1}^e + \gamma^{t_{g_{j+1}}-t_{g_j}} \cdot \max_{g' \in \mathcal{G}} Q^*(s_{t_{g_{j+1}}}, g') \mid s_{t_{g_j}} = s, g_j = g \right]$$

and

$$Q^*(s, a, g) = \mathbb{E} \left[r_t^i + \gamma \cdot \max_{a' \in \mathcal{A}} Q_i^*(s_{t+1}, g_j, a') \mid s_t = s, g_j = g, a_t = a, t \in [t_{g_j}, t_{g_{j+1}}] \right].$$

Here $\gamma \in [0, 1]$ is a discount factor, and the expectations are taken over the randomness of the reward and the state transition,

We use deep neural networks to approximate the two Q-value functions as $Q^*(s, a, g) \approx Q(s, a, g; \theta_i)$ and $Q^*(s, g) \approx Q(s, g; \theta_e)$. The parameters θ_i and θ_e are optimized to minimize the following quadratic loss functions:

$$L_i(\theta_i) = \frac{1}{2|D^i|} \sum_{(s, a, g, s', r^i) \in D^i} [(y^i - Q(s, a, g; \theta_i))^2] \\ y^i = r^i + \gamma \cdot \max_{a' \in \mathcal{A}} Q_i(s', a', g; \theta_i) \quad (3)$$

and

$$L_e(\theta_e) = \frac{1}{2|D^e|} \sum_{(s, g, s', r^e) \in D^e} [(y^e - Q(s, g; \theta_e))^2] \\ y^e = r^e + \gamma \cdot \max_{g' \in \mathcal{G}} Q(s', g'; \theta_e). \quad (4)$$

Here, D^e, D^i are the replay buffers storing dialogue experience for training top-level and low-level policies.

Optimization of parameters θ_i and θ_e can be done by stochastic gradient descent on the two loss functions in Equations (3) and (4). The gradients of the two loss functions w.r.t their parameters are

$$\nabla_{\theta_i} L_i = \frac{1}{|D^i|} \sum_{(s, a, g, s', r^i) \in D^i} \left[\nabla_{\theta_i} Q(s, a, g; \theta_i) \cdot (y^i - Q_i(s, a, g; \theta_i)) \right]$$

and

$$\nabla_{\theta_e} L_e = \frac{1}{|D^e|} \sum_{(s, g, s', r^e) \in D^e} \left[\nabla_{\theta_e} Q(s, g; \theta_e) \cdot (y^e - Q_e(s, g; \theta_e)) \right].$$

To avoid overfitting, we also add L_2 -regularization to the objective functions above.

B Sample Dialogue

Table 3: Sample dialogue by the m-HRL agent interacting with real user: bolded slots are the joint constraints between two subtasks. (*agt*: Agent, *usr*: User)

User Goal
<pre> reserve-hotel subtask: { "request_slots": { "hotel_price": "?", "hotel_date_checkout": "?", "hotel_name": "?", "hotel_amenity_wifi": "?" } "inform_slots": { "hotel_date_checkin": "2016-09-22", "hotel_city": "Curitiba", "hotel_numberofpeople": "4" } } </pre>
<pre> book-flight-ticket subtask: { "request_slots": { "price": "?", "return_time_dep": "?", "return_date_dep": "?", "depart_time_dep": "?", "seat": "?" } "inform_slots": { "or_city": "Lima", "dst_city": "Curitiba", "numberofpeople": "4", "depart_date_dep": "2016-09-22" } } </pre>
m-HRL agent
<pre> agt: Hi, may I know your departure city and destination city? usr: I depart from Lima to Curitiba. agt: What is your departure date? usr: I am leaving on 2016-09-22. agt: There is one flight available at 8:00:00AM. usr: What date is available for return flight? agt: There is one return flight on 2016-09-27. usr: What is the time for the return flight? agt: It is 00:00:00AM. usr: What seat is available? agt: Only Economy seat now. usr: We need 4 tickets for the flight. agt: Total is 1201 dollars for flight. inform(subtask_complete): Subgoal reached! usr: What hotel is available there? agt: There is one hotel called Starlight. usr: Is the wifi available for the hotel? agt: Yes, it is free. usr: What's the price for the hotel? agt: It is 1541 dollars for hotel. usr: Could you please book the trip for me? agt: Okay. Trip tickets are issued. usr: Thanks! agt: Thanks! </pre>