

Strategy and Policy Learning for Non-Task-Oriented Conversational Systems

Zhou Yu, Ziyu Xu, Alan W Black and Alex I. Rudnicky

School of Computer Science

Carnegie Mellon University

{zhouyu, awb, air}@cs.cmu.edu, ziyux@andrew.cmu.edu

Abstract

We propose a set of generic conversational strategies to handle possible system breakdowns in non-task-oriented dialog systems. We also design policies to select these strategies according to dialog context. We combine expert knowledge and the statistical findings derived from data in designing these policies. The policy learned via reinforcement learning outperforms the random selection policy and the locally greedy policy in both simulated and real-world settings. In addition, we propose three metrics for conversation quality evaluation which consider both the local and global quality of the conversation.

1 Introduction

Non-task-oriented conversational systems do not have a stated goal to work towards. Nevertheless, they are useful for many purposes, such as keeping elderly people company and helping second language learners improve conversation and communication skills. More importantly, they can be combined with task-oriented systems to act as a transition smoother or a rapport builder for complex tasks that require user cooperation. There are a variety of methods to generate responses for non-task-oriented systems, such as machine translation (Ritter et al., 2011), retrieval-based response selection (Banchs and Li, 2012), and sequence-to-sequence recurrent neural network (Vinyals and Le, 2015). However, these systems still produce utterances that are incoherent or inappropriate from time to time. To tackle this problem, we propose a set of conversational strategies, such as switching topics, to avoid possible inappropriate responses (breakdowns). After we have a set of strategies, which strategy to perform according to

the conversational context is another critical problem to tackle. In a multi-turn conversation, the user experience will be affected if the same strategy is used repeatedly. We experimented on three policies to control which strategy to use given the context: a random selection policy that randomly selects a policy regardless of the context, a locally greedy policy that focuses on local context, and a reinforcement learning policy that considers conversation quality both locally and globally. The strategies and policies are applicable for non-task-oriented systems in general. The strategies can prevent a possible breakdown, and the probability of possible breakdowns can be calculated using different metrics according to different systems. For example, a neural network generation system (Vinyals and Le, 2015) can use the posterior probability to decide if the generated utterance is possibly causing a breakdown, thus replacing it with a designed strategy. In this paper, we implemented the strategies and policies in a keyword retrieval-based non-task-oriented system. We used the retrieval confidence as the criteria to decide whether a strategy needed to be triggered or not.

Reinforcement learning was introduced to the dialog community two decades ago (Biermann and Long, 1996) and has mainly been used in task-oriented systems (Singh et al., 1999). Researchers have proposed to design dialogue systems in the formalism of Markov decision processes (MDPs) (Levin et al., 1997) or partially observable Markov decision processes (POMDPs) (Williams and Young, 2007). In a stochastic environment, a dialog system's actions are system utterances, and the state is represented by the dialog history. The goal is to design a dialog system that takes actions to maximize some measure of system reward, such as task completion rate or dialog length. The difficulty of such modeling lies in the state representation. Representing the dialog by the entire history is often neither feasible nor

conceptually useful, and the so-called belief state approach is not possible, since we do not even know what features are required to represent the belief state. Previous work (Walker et al., 1998) has largely dealt with this issue by imposing prior limitations on the features used to represent the approximate state. In this paper, instead of focusing on task-oriented systems, we apply reinforcement learning to design a policy to select designed conversation strategies in a non-task-oriented dialog systems. Unlike task-oriented dialog systems, non-task-oriented systems have no specific goal that guides the interaction. Consequently, evaluation metrics that are traditionally used for reward design, such as task completion rate, are no longer appropriate. The state design in reinforcement learning is even more difficult for non-task-oriented systems, as the same conversation would not occur more than once; one slightly different answer would lead to a completely different conversation; moreover there is no clear sense of when such a conversation is “complete”. We simplify the state design by introducing expert knowledge, such as not repeating the same strategy in a row, as well as statistics obtained from conversational data analysis.

We implemented and deployed a non-task-oriented dialog system driven by a statistical policy to avoid possible system breakdowns using designed general conversation strategies. We evaluated the system on the Amazon Mechanical Turk platform with metrics that consider both the local and the global quality of the conversation. In addition, we also published the system source code and the collected conversations ¹.

2 Related Work

Many generic conversational strategies have been proposed in previous work to avoid generating incoherent utterances in non-task-oriented conversations, such as introducing new topics (e.g. “Let’s talk about favorite foods!”) in (Higashinaka et al., 2014), asking the user to explain missing words (e.g. “What is SIGDIAL?”) (Maria Schmidt and Waibel, 2015). In this paper, we propose a set of generic strategies that are inspired by previous work, and test their usability on human users. No researcher has investigated thoroughly on which strategy to use in different conversational contexts. Compared to task-oriented dialog systems, non-

task-oriented systems have more varied conversation history, which are thus harder to formulate as a mathematical problem. In this work, we propose a method to use statistical findings in conversational study to constrain the dialog history space and to use reinforcement learning for statistical policy learning in a non-task-oriented conversation setting.

To date, reinforcement learning is mainly used for learning dialogue policies for slot-filling task-oriented applications such as bus information search (Lee and Eskenazi, 2012), restaurant recommendations (Jurčiček et al., 2012), and sight-seeing recommendations (Misu et al., 2010). Reinforcement learning is also used for some more complex systems, such as learning negotiation policies (Georgila and Traum, 2011) and tutoring (Chi et al., 2011). Reinforcement learning is also used in question-answering systems (Misu et al., 2012). Question-answering systems are very similar to non-task-oriented systems except that they do not consider dialog context in generating responses. They have pre-existing questions that the user is expected to go through, which limits the content space of the dialog. Reinforcement learning has also been applied to a non-task-oriented system for deciding which sub-system to choose to generate a system utterance (Shibata et al., 2014). In this paper, we used reinforcement learning to learn a policy to sequentially decide which conversational strategy to use to avoid possible system breakdowns.

The question of how to evaluate conversational systems has been under discussion throughout the history of dialog system research. Task completion rate is widely used as the conversational metric for task oriented systems (Williams and Young, 2007). However, it is not applicable for non-task-oriented dialog systems which don’t have a task. *Response appropriateness (coherence)* is a widely used manual annotation metric (Yu et al., 2016) for non-task-oriented systems. However, this metric only focuses on the utterance level conversational quality and is not automatically computable. Perplexity of the language model is an automatically computable metric but is hard to interpret (Vinyals and Le, 2015). In this paper, we propose three metrics: *turn-level appropriateness*, *conversational depth* and *information gain*, which access both the local and the global conversation quality of a non-task-oriented conversation. *Information*

¹www.cmuticktock.org

gain is automatically quantifiable. We use supervised machine learning methods to build automatic detectors for *turn level appropriateness* and *conversational depth*. All three of the metrics are general enough to be applied to any non-task-oriented system.

3 Conversational Strategy Design

We implemented ten strategies in total for response generation. The system only selects among Strategy 1-5 if their trigger conditions are met. If more than one strategy is eligible, the system selects the higher ranked strategy. The rank of the strategies, shown in the following list, is determined via expert knowledge. The system only selects among Strategy 6-10 if Strategy 1-5 cannot be selected. This rule reduces the design space of all policies. We design three different versions of the surface form for each strategy, so the user would get a slightly different version every time, thus making the system seem less robotic.

We implemented these strategies in TickTock (Yu et al., 2015). TickTock is a non-task-oriented dialog system that takes typed text as the input and produces text as output. It performs anaphora detection and candidate re-ranking with respect to history similarity to track conversation history. For a detailed system description, please refer to (Yu et al., 2016). This version of TickTock took the form of a web-API, which we put on Amazon Mechanical Turk platform to collect data from a large number of users. The system starts the conversation by proposing a topic to discuss. The topic is randomly selected from five designed topics: movies, music, politics, sports and board games. We track the topic of the conversation throughout the interaction. Each conversation has more than 10 turns. Table 1 is an example conversation of TickTock talking with a human user. We describe the ten strategies with their ranking order in the following.

1. **Match Response** (continue): In a keyword-based system, the retrieval confidence is the weighted score of all the matching keywords from the user input and the chosen utterance from the database. When the retrieval confidence score is higher than a threshold (0.3 in our experiment), we use the retrieved response as the system’s output. If the system is a sequence-to-sequence neural networks system, then we select the output of the system

when the posterior probability of the generated response is higher than a certain threshold.

2. **Don’t Repeat** (no repeat): When users repeat themselves, the system confronts them by saying: “You already said that!”.
3. **Ground on Named Entities** (named entity) A lot of raters assume that TickTock can answer factual questions, so they ask questions such as “Which state is Chicago in?” and “Are you voting for Clinton?”. We use the Wikipedia knowledge base API to tackle such questions. We first perform a shallow parsing to find the named entity in the sentence, and then we search the named entity in a knowledge base, and retrieve the corresponding short description of it. Finally we design several templates to generate sentences using the obtained short description of the named entity. The resulting output can be “Are you talking about the city in Illinois?” and “Are you talking about Bill Clinton, the 42rd president of the United States, or Hillary Clinton, a candidate for the Democratic presidential nomination in the 2016 election?”. This strategy is considered one type of grounding strategy in human conversations. Users feel like they are understood when this strategy is triggered correctly. In addition, we make sure we never ground the same named-entity twice in single conversation.
4. **Ground on Out of Vocabulary Words** (oov) If we find that the user utterance contains a word that is out of our vocabulary, such as “confrontational”. Then TickTock will ask: “What is confrontational?”. We expand our vocabulary with the new user-defined words continuously, so we will not ask for grounding on the same word twice.
5. **React to Single-word Sentence** (short answer) We found that some users type in meaningless single words such as ‘d’, ‘dd’, or equations such as ‘1+2=’. TickTock will reply: “Can you be serious and say things in a complete sentence?” to deal with such condition.
6. **Switch Topic** (switch) TickTock proposes a new topic other than the current topic, such

as “sports” or “music”. For example: “Let’s talk about sports.” If this strategy is executed, we will update the tracked topic to the new topic introduced.

7. **Initiate Activities** (initiation) TickTock invites the user to do an activity together. Each invitation is designed to match the topic of the current conversation. For example, the system would ask: “Do you want to see the latest Star Wars movie together?” when it is talking about movies with a user.
8. **End topics with an open question** (end): TickTock closes the current topic and asks an open question, such as “ Sorry I don’t know. Could you tell me something interesting?”.
9. **Tell A Joke** (joke): TickTock tells a joke such as: “Politicians and diapers have one thing in common. They should both be changed regularly, and for the same reason”. The jokes are designed with respect to different topics as well. The example joke is related to the topic “politics”.
10. **Elicit More Information** (more): TickTock asks the user to say more about the current topic, using utterances such as “ Could we talk more about that?”.

4 Strategy Design

As a baseline policy, we use a random selection policy that randomly chooses among Strategies 6-10 whenever Strategies 1-5 are not applicable. In the conversations collected using the baseline, we found that the sentiment polarity of the utterance has an influence on which strategy to select. People tend to rate the *switch* strategy more favorably if there is negative sentiment in the previous utterances. For example:

TickTock: Hello, I really like politics. Let’s talk about politics.

User: No, I don’t like politics.

TickTock: Why is that?

User: I just don’t like politics.

TickTock: OK, how about we talk about movies?

In another scenario, when all the previous three utterances are positive, the *more* strategy (e.g.

Do you want to talk more about that?) is preferred over the *switch* strategy (e.g. Do you like movies?).

We set out to find the optimum strategy given the context which is the sentiment polarity of the previous three utterances. We found all the scenarios when Strategy 6-10 are triggered, then we generate five different versions of the conversations by replacing the original used strategy with Strategies 6-10. We asked workers on Amazon Mechanical Turk to rate the strategy’s appropriateness given three previous utterances. For each conversation, we collected ratings from three different raters and used the majority vote as the final rating. Then we constructed a table of a distribution of the probability of each strategy with respect to the context. We collected 10 ratings for each strategy under each context. We use the Vader (Hutto and Gilbert, 2014) sentiment predictor for automatic sentiment prediction. The output of the sentiment predictor is a label with three categories: positive (pos), negative (neg) and neutral (neu).

We found that the output of the rating task supports our hypothesis: different strategies are preferred with respect to different sentiment context. In Table 3, we show the distribution of appropriateness ratings for Strategy 6-10 in a context when all the previous utterances are positive. Users rated the *more* strategy more appropriate than the *end* strategy and the *switch* strategy. One interesting observation is that the *joke* strategy is rated poorly. We examined the cases in which it is used and found that the low appropriateness is mostly the result of being unexpected. The *initiation* strategy can be appropriate when the activity fits the previous content semantically. In another sentiment context, when there are consecutive negative utterances, the *switch* strategy and the *end* strategy are preferred. We can see that which strategy to use is heavily dependent on the immediately sentiment context of the conversation. Sentiment polarity captures some conversational level information which is a discriminating factor. We then use these findings to design the locally greedy policy. The system chooses the strategy that is rated as the most appropriate given the context. The context is the sentiment polarity of previous three utterances.

We conducted another Amazon Mechanical Turk study to test if sentiment context beyond three utterances would influence the preferred strategy. In order to reduce the work load, we

Turn	Participant	TickTock
1	I am not sure. What’s on your mind?	Well, you know, I’m not sure. Let me tell you something, politicians and diapers have one thing in common. They should both be changed regularly, and for the same reason.
2	Wow. That’s funny and insightful.	Uhh, I can’t answer this. Can we talk about something else?
3	How about movies. Have you seen any good movies?	Of course. What do you think?
4	I love super hero movies. Have you seen Deadpool?	Are you talking about Deadpool,the animation character?
5	Yes, exactly. They made a movie of him.	That sounds really interesting.

Table 1: An example conversation with TickTock in the real-world setting

Turn	A.L.I.C.E.	TickTock
1	When do I do it?	I don’t know.
2	That is interesting.	Do you want to go to my friend’s party tonight?
3	Sure, just download me to your notebook or PDA and take me along.	Sounds great. Do you want to play board games together some time.
4	I like to play the Turing Game.	I don’t know how to play it. Is it difficult to play?
5	Interesting gossip. It could be.	Possibly.

Table 2: An example conversation of TickTock in the simulated setting

tested one condition: the previous three utterances are all positive. We provide the complete conversation history of that dialog to the raters. We found that strategies used most recently are rated less favorably if used again. This motivates us to include information that relates to the previous strategy usage and a longer history to design policy in the reinforcement learning setting.

Strategy	App	Inter	Inapp
switch	0.1	0.3	0.6
initiation	0.2	0.4	0.4
joke	0.1	0.2	0.7
end	0.1	0.3	0.6
more	0.4	0.5	0.1

Table 3: Appropriateness rating distribution when the recent three utterances are positive.

5 Reinforcement Learning

We model the conversation process as a Markov Decision Process (MDP)-based problem, so we can use reinforcement learning to learn a conversational policy that makes sequential decisions by considering the entire context. We used Q-learning, a model-free method to learn the conversational policy for our non-task-oriented conversational system.

In reinforcement learning, the problem is defined as $(S, A, R, \gamma, \alpha)$, where S is the set of states that represents the system’s environment, in this case the conversational context. A is a set of actions available per state. In our setting, the actions

are strategies available. By performing an action, the agent can move from one state to another. Executing an action in a specific state provides the agent with a reward (a numerical score), $R(s, a)$. The goal of the agent is to maximize its total reward. It does this by learning which action is optimal to take for each state. The action that is optimal for each state is the action that has the highest long-term reward. This reward is a weighted sum of the expected values of the rewards of all future steps starting from the current state, where the discount factor γ is a number between 0 and 1 that trades off the importance of sooner versus later rewards. γ may also be interpreted as the likelihood to succeed (or survive) at every step. The algorithm therefore has a function that calculates the quantity of a state-action combination, $Q : S \times A \rightarrow R$. The core of the algorithm is a simple value iteration update. It assumes the old value and makes a correction based on the new information at each time step, t . See Equation (1) for details of the iteration function.

The critical part of the modeling is to design appropriate states and the corresponding reward function. We reduce the number of the states by incorporating expert knowledge and the statistical findings in our analysis. We used another chatbot, A.L.I.C.E.² as a user simulator in the training process. We include features: turn index, times each strategy was previously executed, and the sentiment polarity of previous three utterances. We constructed the reward table based on the statis-

²<http://alice.pandorabots.com/>

$$Q_{t+1}(s_t, a_t) \leftarrow Q_t(s_t, a_t) + \alpha_t(s_t, a_t) \cdot \left(R_{t+1} + \gamma \max_a Q_t(s_{t+1}, a) - Q_t(s_t, a_t) \right) \quad (1)$$

$$\text{Turn-level appropriateness} * 10 + \text{Conversational depth} * 100 + \text{round}(\text{Information gain}, 5) * 30 \quad (2)$$

tics collected from the previous experiment. In order to make the reward table tractable, we imposed some of the rules we constructed based on expert knowledge. For example, if certain strategies have been used before, then the reward of using it again is reduced. If the trigger condition of Strategy 1-5 is met, the system chooses them over Strategy 6-10. This may result in some less optimum solutions, but reduces the state space and action space considerably. During the training process, we constrained the conversation to be 10 turns. The reward function is only given at the end of the conversation, it is a combination of the automatic predictions of the three metrics that consider the conversation quality both locally and globally, discussed them in detail in the next section. It took 5000 conversations for the algorithm to converge. We looked into the learned Q table and found that the policy prefers the strategy that uses less frequently if the context is fixed.

6 Evaluation Metrics

In the learning process of the reinforcement learning, we use a metric which is a combination of three metrics: *turn-level appropriateness*, *conversational depth* and *information gain*. *Conversational depth* and *information gain* measure the quality of the conversation across multiple turns. Since we use another chatbot as the simulator, making sure the overall conversation quality is accessed is critical. All three metrics are related to each other but cover different aspects of the conversation. We used a weighted score of the three metrics for the learning process, which is shown in Equation (2). The coefficients are chosen based on empirical heuristics. We built automatic predictors for *turn-level appropriateness* and *conversational depth* based on annotated data as well.

6.1 Turn-Level Appropriateness

Turn-level appropriateness reflects the coherence of the system’s response in each conversational turn. See Table 4 for the annotation scheme. The inter-annotator agreement between the two experts is relatively high ($\kappa = 0.73$). We collapse

the “Appropriate” and “Interpretable” labels into one class and formulate the appropriateness detection as a binary classification problem. Our designed policies and strategies intend to avoid system breakdowns (the inappropriate responses), so we built this detector to tell whether a system response is appropriate or not.

We annotated the appropriateness for 1256 turns. We balance the ratings by generating more inappropriate examples by randomly pairing two utterances. In order to reduce the variance of the detector, we use five-fold cross-validation and a Z-score normalizer to scale all the features into the same range. We use early fusion, which simply concatenates all feature vectors. We use a v-Support Vector (Chang and Lin, 2011) with a RBF Kernel to train the detector. The performance of the automatic appropriateness detector is 0.73 in accuracy while the accuracy of the majority vote is 0.5.

We use three sets of features: the strategy used in the response, the word counts of both the user’s and TickTock’s utterances, and the utterance similarity features. The utterance similarity features consist of a feature vector obtained from a word2vec model (Mikolov et al., 2013), the cosine similarity score between the user utterance and the system response, and the similarity scores between the user response and all the previous system responses. For the word2vec model, we trained a 100-dimension model using the collected data.

6.2 Conversational Depth

Conversational depth reflects the number of consecutive utterances that share the same topic. We design an annotation scheme (Table 5) based on the maximum number of consecutive utterances on the same topic. We annotate conversations into three categories: “Shallow”, “Intermediate” and “Deep”. The annotation agreement between the two experts is moderate ($\kappa = 0.45$). Users manually labeled 100 conversations collected using TickTock. We collapse “Shallow” and “Intermediate” into one category and formulate the

Label	Definition	Example
Inappropriate (Inapp)	Not coherent with the user utterance	<i>Participant</i> : How old are you? <i>TickTock</i> : Apple.
Interpretable (Inter)	Related and can be interpreted	<i>Participant</i> : How old are you? <i>TickTock</i> : That’s too big a question for me to answer.
Appropriate (App)	Coherent with the user utterance	<i>Participant</i> : How is the weather today? <i>TickTock</i> : Very good.

Table 4: Appropriateness rating scheme.

Conv. depth	Consecutive utterances
Shallow	< 6
Intermediate	[7, 10]
Deep	> 10

Table 5: Conversational depth annotation scheme

problem as a binary classification problem. We use the same machine learning setting as the turn level appropriateness predictor. The performance of the automatic conversational depth detector has a 72.7% accuracy, while the majority vote baseline accuracy is 63.6%. The conversational depth detector has three types of features:

1. The number of dialogue exchanges between the user and TickTock and the number of times TickTock uses the *continue*, *switch* and *end* strategy.
2. The count of a set of keywords in the conversation. The keywords are “sense”, “something” and interrogative pronouns, such as “when”, “who”, “why”, etc. “Sense” often occurs in sentence, such as “You are not making any sense” and “something” often occurs in sentence, such as “Can we talk about something else?” or “Tell me something you are interested in.”. Both of them indicate a possible change of a topic. Interrogative pronouns are usually involved in questions that probe users to go deep into the current topic.
3. We convert the entire conversation into a vector using doc2vec and also include the cosine similarity scores between adjacent responses of the conversation.

6.3 Information Gain

Information gain reflects the number of unique words that are introduced into the conversation from both the system and the user. We believe

that the more information the conversation has, the better the conversational quality is. This metric is calculated automatically by counting the number of unique words after the utterance is tokenized.

7 Results and Analysis

We evaluate the three policies with respect to three evaluation metrics: *turn-level appropriateness*, *conversational depth* and *information gain*. We show the results in the simulated setting in Table 6 and the real-world setting in Table 7. In the simulated setting, users are simulated using a chatbot, A.L.I.C.E.. We show an example simulated conversion in Table 2. In the real-world setting, the users are people recruited on Amazon Mechanical Turk. We collected 50 conversations for each policy. We compute *turn-level appropriateness* and *conversational depth* using automatic predictors in the simulated setting and use manual annotations in the real-world setting.

The policy learned via reinforcement learning outperforms the other two policies in all three metrics with statistical significance ($p < 0.05$) in both the simulated setting and the real-world setting. The percentage of inappropriate turns decreases when the policy considers context in selecting strategies. However, the percentage of appropriate utterances is not as high as we hoped. This is due to the fact that in some situations, no generic strategy is appropriate. For example, none of the strategies can produce an appropriate response for a content-specific question, such as “What is your favorite part of the movie?” However, the *end* strategy can produce a response, such as: “Sorry, I don’t know, tell me something you are interested.” This strategy is considered “Interpretable” which in turn saves the system from a breakdown. The goal of designing strategies and policies is to avoid system breakdowns, so using the *end* strategy is a good choice in such a situation. These generic strategies are designed to

Policy	Appropriateness	Conversational depth	Info gain
Random Selection	62%	32%	50.2
Locally Greedy	72%	34%	62.4
Reinforcement Learning	82%	45%	68.2

Table 6: Performance of different policies in the simulated setting

Policy	App	Inter	Inapp	Conversational depth	Info gain
Random Selection	30%	36%	32%	30%	56.3
Locally Greedy	30%	42%	27%	52%	71.7
Reinforcement Learning	34%	43%	23%	58%	73.2

Table 7: Performance of different policies in the real-world setting.

avoid system breakdowns, so some times they are not “Appropriate”, but only “Interpretable”.

Both the reinforcement learning policy and the locally greedy policy outperform the random selection policy with a huge margin in conversational depth. The reason is that they take context into consideration in selecting strategies, while the random selection policy uses the *switch* strategy randomly without considering the context. As a result, it cannot keep the user on the same topic for long. However, the reinforcement learning policy only outperforms the locally greedy policy with a small margin. Because there are cases when the user has very little interest in a topic, the reinforcement learning policy will switch the topic to satisfy the *turn-level appropriateness* metric, while the locally greedy policy seldom selects the *switch* strategy according to the learned statistics.

The reinforcement learning policy has the best performance in terms of information gain. We believe the improvement mostly comes from using the *more* strategy appropriately. The *more* strategy elicits more information from the user compared to other strategies in general.

In Table 2, we can see that the simulated user is not as coherent as a human user. In addition, the simulated user is less expressive than a real user, so the depth of the conversation is generally lower in the simulated setting than in the real-world setting.

8 Conclusion and Future Work

We design a set of generic conversational strategies, such as switching topics and grounding on named-entities, to handle possible system breakdowns in any non-task-oriented system. We also learn a policy that considers both the local and global context of the conversation for strategy

selection using reinforcement learning methods. The policy learned by reinforcement learning outperforms the locally greedy policy and the random selection policy with respect to three evaluation metrics: turn-level appropriateness, conversational depth and information gain.

In the future, we wish to consider user’s engagement in designing the strategy selection policy in order to elicit high quality responses from human users.

References

- Rafael E Banchs and Haizhou Li. 2012. Iris: a chat-oriented dialogue system based on the vector space model. In *Proceedings of the ACL 2012 System Demonstrations*, pages 37–42. Association for Computational Linguistics.
- Alan W Biermann and Philip M Long. 1996. The composition of messages in speech-graphics interactive systems. In *Proceedings of the 1996 International Symposium on Spoken Dialogue*, pages 97–100.
- Chih-Chung Chang and Chih-Jen Lin. 2011. LIB-SVM: a library for support vector machines. *ACM Transactions on Intelligent Systems and Technology (TIST)*, 2(3):27.
- Min Chi, Kurt VanLehn, Diane Litman, and Pamela Jordan. 2011. Empirically evaluating the application of reinforcement learning to the induction of effective and adaptive pedagogical strategies. *User Modeling and User-Adapted Interaction*, 21(1-2):137–180.
- Kallirroi Georgila and David R Traum. 2011. Reinforcement learning of argumentation dialogue policies in negotiation. In *INTERSPEECH*, pages 2073–2076.
- Ryuichiro Higashinaka, Kenji Imamura, Toyomi Meguro, Chiaki Miyazaki, Nozomi Kobayashi, Hiroaki

- Sugiyama, Toru Hirano, Toshiro Makino, and Yoshihiro Matsuo. 2014. Towards an open-domain conversational system fully based on natural language processing. In *COLING*, pages 928–939.
- Clayton J Hutto and Eric Gilbert. 2014. Vader: A parsimonious rule-based model for sentiment analysis of social media text. In *Eighth International AAI Conference on Weblogs and Social Media*.
- Filip Jurčiček, Blaise Thomson, and Steve Young. 2012. Reinforcement learning for parameter estimation in statistical spoken dialogue systems. *Computer Speech & Language*, 26(3):168–192.
- Sungjin Lee and Maxine Eskenazi. 2012. Pomdp-based let’s go system for spoken dialog challenge. In *Spoken Language Technology Workshop (SLT), 2012 IEEE*, pages 61–66. IEEE.
- Esther Levin, Roberto Pieraccini, and Wieland Eckert. 1997. Learning dialogue strategies within the markov decision process framework. In *Automatic Speech Recognition and Understanding, 1997. Proceedings., 1997 IEEE Workshop on*, pages 72–79. IEEE.
- Jan Niehues Maria Schmidt and Alex Waibel. 2015. Towards an open-domain social dialog system. In *Proceedings of the 6th International Workshop Series on Spoken Dialog Systems*, pages 124–129.
- Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*.
- Teruhisa Misu, Komei Sugiura, Kiyonori Ohtake, Chiori Hori, Hideki Kashioka, Hisashi Kawai, and Satoshi Nakamura. 2010. Modeling spoken decision making dialogue and optimization of its dialogue strategy. In *Proceedings of the 11th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 221–224. Association for Computational Linguistics.
- Teruhisa Misu, Kallirroi Georgila, Anton Leuski, and David Traum. 2012. Reinforcement learning of question-answering dialogue policies for virtual museum guides. In *Proceedings of the 13th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 84–93. Association for Computational Linguistics.
- Alan Ritter, Colin Cherry, and William B Dolan. 2011. Data-driven response generation in social media. In *Proceedings of the conference on empirical methods in natural language processing*, pages 583–593. Association for Computational Linguistics.
- Tomohide Shibata, Yusuke Egashira, and Sadao Kurohashi. 2014. Chat-like conversational system based on selection of reply generating module with reinforcement learning. In *Proceedings of the 5th International Workshop Series on Spoken Dialog Systems*, pages 124–129.
- Satinder P Singh, Michael J Kearns, Diane J Litman, and Marilyn A Walker. 1999. Reinforcement learning for spoken dialogue systems. In *Nips*, pages 956–962.
- Oriol Vinyals and Quoc Le. 2015. A neural conversational model. *ICML Deep Learning Workshop 2015*.
- Marilyn A Walker, Jeanne C Fromer, and Shrikanth Narayanan. 1998. Learning optimal dialogue strategies: A case study of a spoken dialogue agent for email. In *Proceedings of the 36th Annual Meeting of the Association for Computational Linguistics and 17th International Conference on Computational Linguistics-Volume 2*, pages 1345–1351. Association for Computational Linguistics.
- Jason D Williams and Steve Young. 2007. Partially observable markov decision processes for spoken dialog systems. *Computer Speech & Language*, 21(2):393–422.
- Zhou Yu, Alexandros Papangelis, and Alexander Rudnicky. 2015. TickTock: A non-goal-oriented multimodal dialog system with engagement awareness. In *Proceedings of the AAI Spring Symposium*.
- Zhou Yu, Ziyu Xu, Alan Black, and Alexander Rudnicky. 2016. Chatbot evaluation and database expansion via crowdsourcing. In *Proceedings of the chatbot workshop of LREC*.