# Goal-oriented dialog systems and Memory : an overview

**Léon-Paul Schaub**[*][†]**, Cindel Vaudapiviz** [*]

[*] LIMSI, CNRS, Université Paris-Saclay, F-91405 Orsay, France
[*] AKIO, 43, rue de Dunkerque Paris 75010
[†]EA 2520 ERTIM, Inalco, 2 rue de Lille 75007 Paris, France
schaub@limsi.fr    cyndel.vaudapiviz@gmail.com

**Résumé**

After recalling generally accepted models for representing human dialogues and the different types of human memory, we review goal oriented dialog systems and chatbots, detailing both dialog models and memory-based systems. From some recent investigations about core human language working hinting at reflective aspects, we explain how we think that the recent advances in adversarial learning could provide an interesting research avenue for improving goal oriented chatbots architecture, by injecting in the usual representation triple: dialog turn, dialogue history and knowledge base, more cognitive oriented aspects inspired by the human memory model.

## 1.   Dialog and memory

**Dialog** acts can be defined as the meaningful exchange of utterance between two or more people where they both take alternatively the role of hearer and of speaker. (Mann et al., 1977) propose the "dialog game model" based on the respective goals of each interlocutor, in which every utterance is expected a priori to contribute to the fulfillment of the (final) dialog goal in addition to providing clues about the personal goal of each interlocutor in turn (Bunt, 2011). Any autonomous systems whose behaviour is more elaborate than a mere set of reflex actions in reaction to changes in his environment, needs to have some sort of dynamic memory in order to take context dependent decisions. In a conversation, the memory is at the center of the decision-making process both for language processing and co-construction of the dialog (Vollmer et al., 2014). A dialogue between two persons has three main drives : empathy, experience and knowledge. Empathy allows a hearer to understand and interpret the speaker emotional state and "normalize" the speaker utterance accordingly. The experience helps the hearer to enrich the context of the current conversation with information from previous interactions (Asghar et al., 2017). Finally, the knowledge is a store of information representing the beliefs one has about the world and human culture. If someone says "Hi", while looking at you, your knowledge informs you that this person is greeting you. The question is : how this trinity works, and how are we capable of processing all this ?

**Memory** is a key element in the process. (Atkinson and Shiffrin, 1968) define memory as "the ability of an intelligent system to record, preserve, and recall past experiences to interpret present experiences." They distinguish three kinds of memory : *(I) sensory* (Klatzky, 1980), the instant perception of a element, *(II) short-term* (Sperling, 1967) and *(III) long-term* (Rudner and Rönnberg, 2008). (Squire and Zola, 1996) precise the model and explains that long-term memory is divided into *(III.a) declaratory* and *(III.b) non-declaratory memory*. (Greenberg and Verfaellie, 2010) refers to *(III.a.1) episodic memory* (autobiographical and egocentric records of past events) and *(III.a.2) semantic memory* (general exocentric records), which together make a sort of self-made knowledge base about the world. (Ebbinghaus, 2013) describes the backup concept that allows information to be stored in the long-term memory. (Miller, 1962) calls non-declaratory memory *(III.b) procedural memory*, which corresponds to what humans unconsciously retain, in particular proprioceptive routines, e.g. when playing a musical instrument. Short-term is actually a generic term that includes (cf Fig. 1) : *(II.a) working memory* (Baddeley, 2010) to process short-term information, the phonological loop (Baddeley et al., 1998) to receive sounds and interpret them and the visio-spatial notebook (Baddeley, 1995) to make a mental representation of the information. Short-term memory lasts only a few seconds (7 to 12) according to experiments [1]. Concepts are transformed into information that is encoded and stored into long-term memory. These are the retention and recall mechanisms (Bodner and Lindsay, 2003). (Bangerter, 2004) talks about joint attention to explain the focus of all interlocutors on the same conversational entity. This attention is focused on the words exchanged in a dialogue. (Cintrón-Valentín and Ellis, 2016) describe it as linguistic salience where a hearer summarizes the main theme of a locutor utterance by selecting the most salient words. During a conversation, when one hears an utterance from the interlocutor : the sensory memory intercepts the most salient information and stores it in only a few miliseconds, before sending it to the working memory. There, it will be preserved until an interpretation can be produced. The working memory acts like a buffer while the long-term memory is interrogated to extract the information needed for interpreting the last utterance. First the episodic memory is searched for similar events. When one is found, the semantic memory is searched for the corresponding general meanings and cultural knowledge. All these information are synthesized back into the working memory, which evaluates what information to keep. This back and forth transfer between working and long term memory stops when enough confidence in the current understanding is achieved and the ensuing action can be decided. The dialogue model of (Clark and Marshall, 1981) has five parts : (1) environment, (2) speech

---

1. https://www.cognifit.com/science/
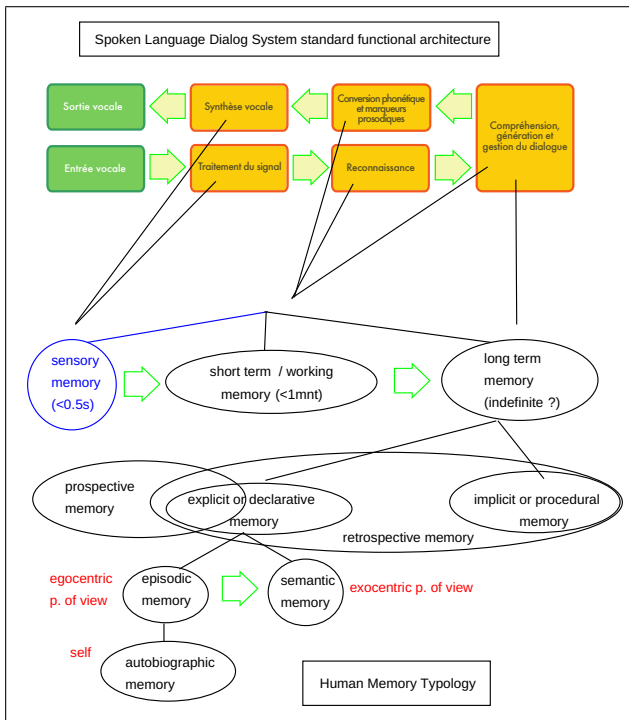cognitive-skills/shortterm-memory

FIGURE 1 – Cognitive memory and spoken dialogue system standard functional architecture
from (Mariani et al., 2012) and www.human-memory.net

turn, (3) history of the dialogue, (4) experience of the interlocutors and (5) knowledge of the interlocutors. If we compare this to the human memory model, we find the following mapping : sensory memory (speech turn), working memory (history of the dialogue), episodic memory (experience), and both semantic and procedural memory (knowledge). The sensory memory is the capacity of the agent to have a representation of the user profile model like in (Janarthanam and Lemon, 2014). It also gets the most salient data in the input and transfers all the information to the working memory. The working memory is the current speaking turn retention and a specific number of previous turns tracks, according to the sensory memory's information. It is the central computer of the information processing during a conversation (Sauseng et al., 2005). Then, the working memory encodes the information it kept in the long-term memory, where are stocked previous conversations (episodic memory) and world knowledge (semantic memory)(Sieber and Krenn, 2010a), with all the user's intentions and the utterances associated, so it can decide the best answer to provide. Then, the working memory retrieves the information encoded in the long-term memory so it can send the information to the response generation processor. The next question now is : what about memory in computational dialog systems ?

## 2. Goal-oriented dialog systems

Compared to chat-oriented dialog system (chatbot), the purpose of which is to maintain a small-talk conversation with a human user, a goal-oriented (GO) dialog system is meant to resolve a pragmatic real-life problem (tickets booking, CRM issue...). In any GO system, there are three elements : the NLU module (understand the user's input), the dialogue policy, including the dialogue management (to decide where to look for the right answer according toNLU module input's comprehension) and the NLG module (how to generate a proper natural language (NL) answer according to what the DM found out). (Wang, 2018) proposes a survey of GO dialogs since ELIZA (Weizenbaum, 1983), the very first dialogue system, until recent systems. The GO motivation is mostly industrial because it is about solving simple (but not always easy) problems of people without having to interact with a human agent. Therefore, we part from the hypothesis that the user is "relatively" collaborative and prioritize the reaching of his.er goal over anything else.

In the early 20', the retrieval-based approach was the most advanced dialog model. A system implementing the model is capable of giving a right answer in a very high number of cases by parsing the question and searching in its knowledge base what is the closest possible answer. It is popular because adapted for reusing information extracted from previous dialogs to compute new replies, like in the emblematic Watson system (Ferrucci et al., 2010). Nevertheless, in a recent experiment (Schaub, 2017), we have observed that better performance than IBM Watson can be obtained with freely available toolkits off the shelf, (sklearn, keras) when evaluated by human on the same task and the same corpus in an industrial context, proving that private solutions don't always work so well. (Young, 2006) show that the POMDP is a good way to model a dialogue system because of the evaluation function inherent to the process, good for accuracy of decision management estimation. (Young et al., 2013) uses it to optimize the agent's choice when providing the answer. Because of it's probabilistic nature, the POMDP fits well with goal oriented dialog requirements : the dialog manager can decide in a non deterministic way the best answer to provide.

Recently, (Wang, 2018) reports that the best performing neuronal approaches are the ones based on LSTM and the ones using reinforcement learning (RL) like (Weisz et al., 2018). Accordingly, the survey from (Chen et al., 2017) identifies LSTM and RL as most performing approaches, in particular in conjunction with active learning (AL) approaches (Asghar et al., 2016), where the learning material provided by the users feedback is sorted and filtered according to the amount of new information it contributes : Fig. 1. The LIHLITH project (Agirre et al., 2018)studies how to modify the dialog system architecture to use lifelong learning for training a chatbot even in production mode.

### 2.1. Dialog models overview

The most common dialog model is (*Dialogue State Tracking* or *belief tracking*). StateNet gets the best results for now (Ren et al., 2018). It is a universal state tracker. (Mrksic et al., 2016) build the first neural system that achieves better results than symbolic ones. The frame tracker model (Schulz et al., 2017) allows the system to free itself from the immediate user intent search. (Ultes et al., 2019) developed the conversational entity model, which abandons traditional dialogue turn to focus only on the re-
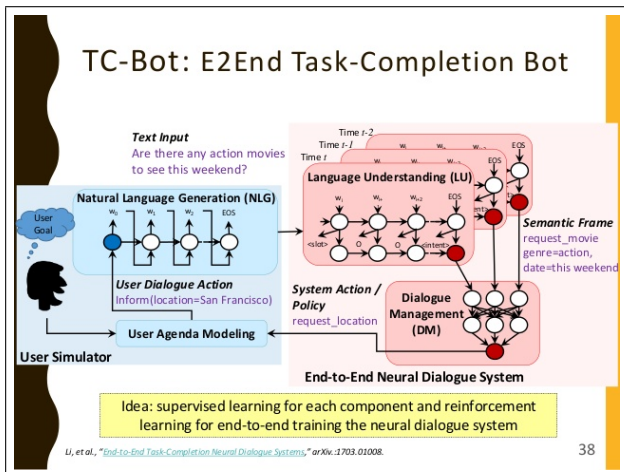
FIGURE 2 – "end-to-end" system model

lations between conversational entities which they define as : "a virtual entity that exists in the context of the conversation and is either a conversational object or a conversational relationship". This is a huge advance in the dialog systems investigation because it would mean that a dialog system must not consider a conversation as a linear list of dialog turns, but as a graph where different slots containing information and relation between information are the "memory" of the conversational agent.

(Asghar et al., 2017) explored emotion in dialogue and showed that analyzing user emotions improved system performance by altering the expected answer according to emotional state of the user. We will introduce works focused on the place of the memory in the conversational agent.

### 2.2. Memory-based dialog systems

The question of a cognitive memory inside an agent has already been raised by (Rhodes, 1997) to specify a prototype agent capable of remembering its interactions, arguing that it would improve the answer decision of the conversational agent. (Sieber and Krenn, 2010b) implement the difference between episodic and semantic memory in a dialogue system. The results show that by building a graph of previous conversations (experience) for the agent, independently of the knowledge base, the agent is able to resolve partially the coreference issue and to understand more accurately the user intentions because it can make the difference between dialog memory and knowledge memory. (Vetulani, 2005) talks about the problem of exception in dialogue : i.e. the ability of a computational system to produce local assumptions from previous speech turns. Some dialogue-specific networks such as "memory networks" (Wan et al., 2018) also show good performances. Recently, (Kim et al., 2019) have combined these memory arrays with Bi-LSTMs (Schuster and Paliwal, 1997) to build an end-to-end system with very satisfactory results when evaluated at the DSTC6 (Hori and Hori, 2017). On this occasion, they showed that building an end-to-end system, where the three dialogue module are merged together can achieve better results than the sepa-

ration of the three. (Zhang et al., 2018) also talk about long-term memory accessible through the content in their system. By making this memory more complex, the agent is able to better distinguish the differences between two records and refine his response accordingly. (Chen et al., 2018) have built a model based on hierarchical memory arrays that allows a conversational agent to make the link between speech turn and long-term memory through a continuous process of knowledge base querying during speech act perception is closer to human behaviour and gives very good results. (El Asri et al., 2017) study in detail the role of cognitive memory in conversational interaction and created a corpus with annotations showing for each speech turn the type of focus. The corpus shows that according to some context, similar inputs can be interpreted and answered differently. We have just seen that the attention in goal oriented chatbot development has been put on the implementation of memory model with cognitive aspects, but this works did not take into account the reflective dimension of language processing hypothesized by some recent works, which we will introduce in the next section in support of our interest for adversarial learning.

### 3. Adversarial learning

Among the three kinds of theories which address the nature of speech units, recent works have put focus on the perceptuo-motor theory (Schwartz et al., 2008), for instance with the COSMO model (Barnaud et al., 2017) were five probabilistic variables define the model of speech units, incorporating both the role of speaker and listener in a coherent coupling. Such works follow the actual trend in favor of sensori-motor models of perception (Gordon et al., 2011) in general. It is interesting to notice that this idea of joining both speaker and listener roles together represents in a sense an instance of adversarial learning, at least for all the times in a dialog, when speech turns overlap, a sign that the interlocutors are "fighting for the floor" (Martine Adda-Decker and Habert, 2008).

Adversarial learning (where two algorithm compete with opposite objectives in a common task) (Borodin and El-Yaniv, 1998) or adversarial neural networks (Goodfellow et al., 2014) have already been tested recently for dialog generation, with interesting performances when compared to reference scores (Liu and Lane, 2018).

We have the intuition that it would be worthwhile to deploy the active learning paradigm in a dialog system at the interface between natural language understanding and dialog management. First because adversarial learning would provide better interpretation, through more focused exploration of the valid interpretation search space. In addition, adversarial learning can often palliate the lack of training data by generating learning data for the adversary by transforming positive examples from the existing learning material into negative training examples for the adversary functionality. And thirdly to explore the possible advantages one could get from having a more cognitive oriented architecture, with reflective aspect that would enable finer distinctions between episodic and semantic memory content. Of course, such experiments will have to guard against the usual risk of divergence from the original data model and

overfitting.

On the practical side, we have identified the recently released RASA STACK (Bocklisch et al., 2017) dialog framework as a good candidate for being used as our experience platform, since it offers convenient support for the three standard dialog functions : understanding, dialog management and generation with convenient modules and interface for state-of-the-art machine learning deployment (POMDP, word embeddings, neural computations) in Python. Furthermore, the framework has an interactive training system, with which the developer can interact to improve decision making.

## 4. Conclusion

After recalling what is a dialog, we have reviewed what we know about the various types of human memory and how they map onto the standard functions found in spoken language systems or chatbots. Then we have presented the state-of-the-art in goal oriented dialog systems, with a focus on memory based approaches. Using as support the recent interest for sensory-motor models of perception, in particular for explaining the nature of speech unit, we have laid out a plan for testing [2] the benefit of having chatbots with a more oriented cognitive model of memory and the adversarial learning paradigm deployed at the interface between the language understanding and dialog management functions.

## 5. References

Agirre, Eneko, Sarah Marchand, Sophie Rosset, Anselmo Peñas, and Mark Cieliebak, 2018. LIHLITH : improving communication skills of robots through lifelong learning. *ERCIM News*, 2018(114).

Asghar, Nabiha, Pascal Poupart, Jesse Hoey, Xin Jiang, and Lili Mou, 2017. Affective neural response generation. *CoRR*, abs/1709.03968.

Asghar, Nabiha, Pascal Poupart, Xin Jiang, and Hang Li, 2016. Online sequence-to-sequence active learning for open-domain dialogue generation. *CoRR*, abs/1612.03929.

Atkinson, R.C. and R.M. Shiffrin, 1968. Human memory. volume 2 of *Psychology of Learning and Motivation*. Academic Press, pages 89 – 195.

Baddeley, Alan, 2010. Working memory. *Current Biology*, 20(4) :R136 – R140.

Baddeley, Alan, Susan Gathercole, and Costanza Papagno, 1998. The phonological loop as a language learning device. *Psychological review*, 105 :158–73.

Baddeley, Alan D., 1995. chapter The psychology of memory. Oxford, England : John Wiley and Sons, pages 3–25.

Bangerter, Adrian, 2004. Using pointing and describing to achieve joint focus of attention in dialogue. *Psychological science*, 15 6 :415–9.

Barnaud, Marie-Lou, Julien Diard, Pierre Bessière, and Jean-Luc Schwartz, 2017. Perceptuo-motor speech units in the brain with cosmo. In *11th International Seminar on Speech Production (ISSP 2017)*.

Bocklisch, Tom, Joey Faulkner, Nick Pawlowski, and Alan Nichol, 2017. Rasa : Open source language understanding and dialogue management. *CoRR*, abs/1712.05181.

Bodner, Glen E. and D. Stephen Lindsay, 2003. Remembering and knowing in context. *Journal of Memory and Language*, 48(3) :563 – 580.

Borodin, Allan and Ran El-Yaniv, 1998. Online computation and competitive analysis.

Bunt, Harry, 2011. The semantics of dialogue acts. In *Proceedings of the Ninth International Conference on Computational Semantics*, IWCS '11. Stroudsburg, PA, USA : ACL.

Chen, Hongshen, Xiaorui Liu, Dawei Yin, and Jiliang Tang, 2017. A survey on dialogue systems : Recent advances and new frontiers. *CoRR*, abs/1711.01731.

Chen, Hongshen, Zhaochun Ren, Jiliang Tang, Yihong Eric Zhao, and Dawei Yin, 2018. Hierarchical variational memory network for dialogue generation. WWW '18. International World Wide Web Conferences Steering Committee.

Cintrón-Valentín, Myrna C. and Nick C. Ellis, 2016. Salience in second language acquisition : Physical form, learner attention, and instructional focus. *Front Psychol*, 7 :1284–1284. 27621715[pmid].

Clark, Herbert H. and Catherine R. Marshall, 1981. Definite knowledge and mutual knowledge. In Aravind K. Joshi, Bonnie L. Webber, and Ivan A. Sag (eds.), *Elements of Discourse Understanding*. Cambridge, UK : Cambridge University Press, pages 10–63.

Ebbinghaus, Hermann, 2013. Memory : a contribution to experimental psychology. *Ann Neurosci*, 20(4) :155–156. 25206041[pmid].

El Asri, Layla, 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. ACL.

Ferrucci, David, Eric Brown, Jennifer Chu-Carroll, James Fan, David Gondek, Aditya A. Kalyanpur, Adam Lally, J. William Murdock, Eric Nyberg, John Prager, Nico Schlaefer, and Chris Welty, 2010. Building watson : An overview of the deepqa project. *AI magazines*, 31(3) :59–79.

Goodfellow, Ian, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio, 2014. Generative adversarial nets. In *Advances in Neural Information Processing Systems 27*. Curran Associates, Inc., pages 2672–2680.

Gordon, G, DM Kaplan, B Lankow, DY Little, J. Sherwin, BA Suter, and L. Thaler, 2011. Toward an integrated approach to perception and action. *Frontiers in Systems Neuroscience*, 5(20).

Greenberg, Daniel L. and Mieke Verfaellie, 2010. Interdependence of episodic and semantic memory : evidence from neuropsychology. *J Int Neuropsychol Soc*, 16(5) :748–753. 20561378[pmid].

Hori, Chiori and Takaaki Hori, 2017. End-to-end

conversation modeling track in DSTC6. *CoRR*, abs/1706.07440.

Janarthanam, Srinivasan and Oliver Lemon, 2014. Adaptive generation in dialogue systems using dynamic user modeling. *Computational Linguistics*, 40(4) :883–920.

Kim, Byoungjae, KyungTae Chung, Jeongpil Lee, Jungyun Seo, and Myoung-Wan Koo, 2019. A bi-lstm memory network for end-to-end goal-oriented dialog learning. *Computer Speech and Language*, 53 :217 – 230.

Klatzky, Roberta, 1980. Human memory : Structures and processes. *The American Journal of Psychology*, 93.

Liu, Bing and Ian Lane, 2018. Adversarial learning of task-oriented neural dialog models. In *Proceedings of the 19th Annual SIGdial Meeting on Discourse and Dialogue*. ACL.

Mann, William C., James A. Moore, and James A. Levin, 1977. A comprehension model for human dialogue. In *Proceedings of the 5th International Joint Conference on Artificial Intelligence - Volume 1*, IJCAI'77. Morgan Kaufmann Publishers Inc.

Mariani, Joseph, Patrick Paroubek, Gil Francopoulo, Aurélien Max, François Yvon, and Pierre Zweigenbaum, 2012. *La langue française à l' Ère du numérique – The French Language in the Digital Age*. Springer.

Martine Adda-Decker, Gilles Adda Patrick Paroubek Philippe Boula de Mareuil, Claude Barras and Benoit Habert, 2008. Annotation and analysis of overlapping speech in political interviews. Marrakech, Morocco : European Language Resources Association (ELRA).

Miller, Brenda, 1962. "physiologie de l'hippocampe". Colloques internationaux du Centre national de la recherche scientifique.

Mrksic, Nikola, Diarmuid Ó Séaghdha, Tsung-Hsien Wen, Blaise Thomson, and Steve J. Young, 2016. Neural belief tracker : Data-driven dialogue state tracking. *CoRR*, abs/1606.03777.

Ren, Liliang, Kaige Xie, Lu Chen, and Kai Yu, 2018. Towards universal dialogue state tracking. *CoRR*, abs/1810.09587.

Rhodes, B. J., 1997. The wearable remembrance agent : a system for augmented memory. In *Digest of Papers. First International Symposium on Wearable Computers*.

Rudner, Mary and Jerker Rönnberg, 2008. The role of the episodic buffer in working memory for language processing. *Cognitive processing*, 9 :19–28.

Sauseng, Paul, Wolfgang Klimesch, Manuel Schabus, and Michael Doppelmayr, 2005. executive functions of working memory. *International Journal of Psychophysiology*, 57(2) :97 – 103. EEG Coherence.

Schaub, Léon-Paul, 2017. *Récupération d'information dans un système de Question-Réponse*. Master's thesis, Inalco, Paris, France.

Schulz, Hannes, Jeremie Zumer, Layla El Asri, and Shikhar Sharma, 2017. A frame tracking model for memory-enhanced dialogue systems. *CoRR*, abs/1706.01690.

Schuster, M. and K.K. Paliwal, 1997. Bidirectional recurrent neural networks. *Trans. Sig. Proc.*, 45(11) :2673–2681.

Schwartz, Jean-Luc, Marc Sato, and Luciano Fadiga, 2008. The common language of speech perception and action. *Revue française de linguistique appliquée*, 13(2) :9–22.

Sieber, Gregor and Brigitte Krenn, 2010a. Episodic memory for companion dialogue. In *Proceedings of the 2010 Workshop on Companionable Dialogue Systems*. ACL.

Sieber, Gregor and Brigitte Krenn, 2010b. Towards an episodic memory for companion dialogue. In Jan Allbeck, Norman Badler, Timothy Bickmore, Catherine Pelachaud, and Alla Safonova (eds.), *Intelligent Virtual Agents*. Berlin, Heidelberg : Springer Berlin Heidelberg.

Sperling, G., 1967. Successive approximations to a model for short term memory. *Acta Psychologica*, 27 :285 – 292.

Squire, Larry and Stuart Zola, 1996. Structure and function of declarative and nondeclarative memory. *Proceedings of the National Academy of Sciences*, 93.

Ultes, Stefan, Pawel Budzianowski, Inigo Casanueva, Lina Maria Rojas-Barahona, Bo-Hsiang Tseng, Yen-Chen Wu, Steve J. Young, and Milica Gasic, 2019. Addressing objects and their relations : The conversational entity dialogue model. *CoRR*, abs/1901.01466.

Vetulani, Zygmunt, 2005. Dialogue processing memory for incident solving in man-machine dialogue. In Leonard Bolc, Zbigniew Michalewicz, and Toyoaki Nishida (eds.), *Intelligent Media Technology for Communicative Intelligence*. Springer Berlin Heidelberg.

Vollmer, Anna Lisa, Jonathan Grizou, Manuel Lopes, Katharina Rohlfing, and Pierre-Yves Oudeyer, 2014. Studying the co-construction of interaction protocols in collaborative tasks with humans. In *The Fourth Joint IEEE International Conference on Development and Learning and on Epigenetic Robotics*.

Wan, Yao, Wenqiang Yan, Jianwei Gao, Zhou Zhao, Jian Wu, and Philip S. Yu, 2018. Improved dynamic memory network for dialogue act classification with adversarial training. *CoRR*, abs/1811.05021.

Wang, Haixun, 2018. An annotated reading list of conversational ai. *Medium*.

Weisz, Gellert, Pawel Budzianowski, Pei-Hao Su, and Milica Gasic, 2018. Sample efficient deep reinforcement learning for dialogue systems with large action spaces. *CoRR*, abs/1802.03753.

Weizenbaum, Joseph, 1983. ELIZA - A computer program for the study of natural language communication between man and machine (reprint). *Commun. ACM*, 26(1) :23–28.

Young, S., 2006. Using pomdps for dialog management. In *2006 IEEE Spoken Language Technology Workshop*.

Young, Steve, Catherine Breslin, Milica Gasic, Matthew Henderson, Dongho Kim, Martin Szummer, Blaise Thomson, Pirros Tsiakoulis, and Eli Tzirkel-Hancock, 2013. Evaluation of statistical pomdp-based dialogue systems in noisy environments.

Zhang, Zheng, Minlie Huang, Zhongzhou Zhao, Feng Ji, Haiqing Chen, and Xiaoyan Zhu, 2018. Memory-augmented dialogue management for task-oriented dialogue systems. *CoRR*, abs/1805.00150.