# Injecting Designers' Knowledge in Conversational Neural Network Systems

Giancarlo A. Xompero[1], Cristina Giannone[1], Fabio Massimo
Zanzotto[2][0000−0002−7301−3596], Andrea Favalli[1], and Raniero Romagnoli[1]

[1] Language Technology Lab, Almawave srl, Rome, Italy
`[first name initial].[last name]@almawave.it`
[2] ART Group, University of Rome Tor Vergata, Rome, Italy
`fabio.massimo.zanzotto@uniroma2.it`

**Abstract.** Sequence-to-sequence neural networks are redesigning dialog managers for Conversational AI in industries. However, industrial applications impose two important constraints: training data are often scarce and the behavior of dialog managers should be strictly controlled and certified. In this paper, we propose the Conversational Logic Injected Neural Network (CLINN). This novel network merges dialog managers "programmed" using logical rules and a Sequence-to-Sequence Neural Network. We experimented with the *Restaurant topic* of the MultiWOZ dataset. Results show that injected rules are effective when training data set are scarce as well as when more data are available.[3]

## 1 Introduction

Sequence-to-sequence neural networks are giving an unprecedented boost to dialog systems and to the adoption of Conversational AI in industries. Sequence-to-sequence dialog systems based on Recurrent Neural Networks (RNNs) have been used to train open domain [9, 7] as well as task-oriented [12] dialog systems. These RNN-based dialog systems have reached interesting results given a sufficiently big set of training data. Transformer-based systems, instead, are less demanding as these can be pre-trained on large datasets and, then, adapted to carry out specific task-oriented dialogs [4, 10, 2]. Due to its interesting performance, Conversational AI is becoming an integral part of business practice across industries[4]. More and more companies are adopting the advantages dialog systems or chatbots bring to customer service, sales as well as workplace assistant.

However, the adoption of conversational AI in industries impose two important constraints on the design of dialog systems: (1) the scarcity of training data and (2) the need for an extreme control on the behavior of dialog systems. In fact, in industrial applications, the scarcity and, sometimes, the complete absence of pre-existing conversation data is the norm. Generally, the Wizard-of-Oz approach for data collecting [11] is adopted to generate training data. This is an expensive process and it is generally

---

[4] https://www.gartner.com/smarterwithgartner/chatbots-will-appeal-to-modern-workers/

not able to provide high quality datasets [8]. On the other hand, the need for an extreme control of dialog systems is generally solved by using dialog systems that can be "programmed" with explicit rules. Undoubtedly, these dialog systems offer extremely precise dialog control in business scenario need and, at the same time, guarantying a satisfying experience for users in covered cases. In this context, design conversational experience is done by defining rules depending on the dialog context and on interpretations of user inputs [6]. Hand-crafted rules ensure generally more control in the conversation flow but do not guarantee scalabily and the generalization given by learning approaches. If dialog interactions are not explicitly modeled, the interaction may miserably fail.

In this paper, we propose to empowering Seq-to-Seq Neural Networks with Conversational Logic Instructions, to satisfy the two industrial constraints on these sequence-to-sequence dialog systems. We adopt a neural dialog manager, based on the Domain Aware Multi-Decoder network [14], adding to it explicit conversational logic instructions to keep human-in-the-loop [13]. The Conversational Logical Injection in Neural Network (CLINN) system combines the generalized power of neural architectures with the control on specific conversational patterns defined by the designers. We experimented with the *Restaurant topic* of the MultiWOZ dataset [1]. We used two different sets of dialogs to allow conversational designers to generate explicit rules. Results show that rules injected are effective in the situation when training data are scarce and, moreover, the defined behaviors on specific conversational patterns are preserved.

## 2   Method and System

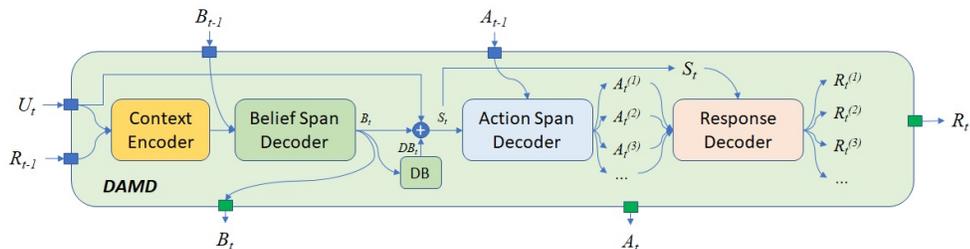### 2.1   Domain Aware Multi-Decoder (DAMD) network



**Fig. 1.** Architecture of the Domain Aware Multi-Decoder (DAMD) network

In this study, we use an end-to-end dialog architecture that includes the concept of *belief span* [5]. The belief span is a sequence of symbols that expresses the belief state at each turn of the dialog. In particular, we rely on the pipeline realized by Zhang et al. [14] that consists of four seq-to-seq modules plus the access to an external database (Fig. 1). The pipeline is applied for each turn of the dialog. It, globally, takes four inputs

$(U_t, R_{t-1}, B_{t-1}, A_{t-1})$ and produces three outputs $(R_t, B_t, A_t)$ where $t$ is the actual turn, $U_t$ is the user utterance, $R_{t-1}$ and $R_t$ are the previous and the current system responses, $B_{t-1}$ and $B_t$ are the previous and the current belief state spans, $A_{t-1}$ and $A_t$ are the previous and the produced system actions. The four modules behave as follows. The *context encoder* encodes the context of the turn $(U_t, R_{t-1})$ in a context vector $c_t$. The *belief span decoder* decodes the previous belief span $B_{t-1}$ and, along with the context vector $c_t$ produces the belief span $B_t$ of current turn. This $B_t$ is used to query the database $DB$ and the answer $DB_t$ is concatenated with $B_t$ to form the internal state $S_t$ of the turn. Then, the *action span decoder* produces the current action $A_t^{(i)}$ by taking into consideration the current state $S_t$ and the previous action $A_{t-1}$. Finally, the *response decoder* emits the final response $R_t^i$ taking into consideration the current state $S_t$ and the corresponding action $A_t^{(i)}$. In [14], multiple actions and multiple responses are produced to increase variability in dialogues and, for this reason, the framework is called multi-action data augmentation.

## 2.2 Injecting Hand-Crafted Knowledge in DAMD

DAMD network offers a tremendous opportunity to inject external knowledge. In fact, the *belief span decoder* transforms the internal context vector $c_t$ and an explicit symbolic previous belief span $B_{t-1}$ in an explicit belief span $B_t$. In the same way, the action span decoder takes in input an explicit, symbolic previous action $A_{t-1}$. As $B_{t-1}$ and $A_{t-1}$ are explicit, these can be easily controlled by an external, symbolic module.
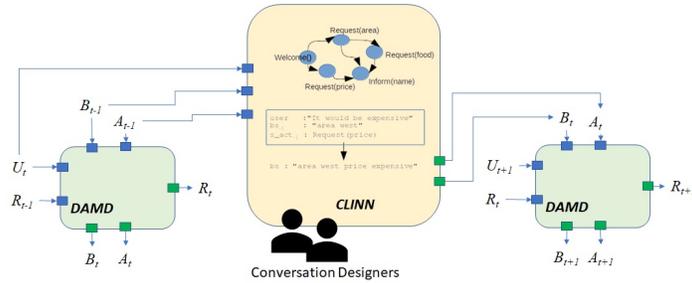


**Fig. 2.** Injecting External Knowledge in DAMD with CLINN

We then propose an external knowledge injector module, that is, our Conversational Logical Injection in Neural Network (CLINN), that allows conversational designers to control the dialog flow with symbolic rules. CLINN acts in between turns, that is, it takes the output and the input of the DAMD network at a given turn $t$ and gives an input to the next step (Fig. 2). CLINN aims to control the next belief state $B_t$ and the action $A_t$ given the previous belief state $B_{t-1}$, the previous action $A_{t-1}$ and the current user utterance $U_t$.

We integrated the CLINN approach into a rule based dialog management system [3]. The rules are derived from the state machine diagram designed by the conversational

designers when they defined the interaction experience in term of tasks and behaviors of the conversational agent. Within the diagram the conversation is defined in term of system actions (i.e. the states) and user input and belief span (in the edges), i.e. the preconditions for changing the state. These are a convenient way for designers to express the conversation behavior they want to mould[5]. In our setting, these diagrams become logical rules that fire when preconditions are matched in the conversation turn. Designing the behaviors for all the possible interactions is very hard and unfruitful. Then, training a neural network can be the solution. However, training a neural network requires a lot of data. Writing symbolic rules is way to inject knowledge in CLINN to boost neural network learning.

## 3 Experiments

### 3.1 Experimental Set-Up

We evaluated CLINN on the MultiWOZ dataset [1] as in Zhang et al. [14]. This dataset is widely used and it has been designed as a human-human task-oriented dialog dataset collected via the Wizard-of-Oz framework. One participant plays the role of the system. The dataset contains conversations on several domains in the area of touristic information (hotel, train, restaurant, taxi,...). Each domain has a set of dialog acts in addition to some general acts such as *greeting* or *goodbye*. Users' and system's interactions are described in term of these dialog acts.

We focused on the restaurant domain of the MultiWOZ dataset that consists of 1200 dialogs for the training set, 61 dialogs for the testing set and 50 dialogs for the validation set. We used two different settings for the training set: (1) a small set of 150 randomly selected dialogs; (2) the full set of 1200 dialogs. These two settings are relevant to study the behavior of our system with few training examples.

In order to simulate the delivering in production environment of a conversational agent, we modeled a state transition diagram, which describes the expected conversational behavior of the agent. The diagram is defined observing some conversational examples in the training set. For the evaluation we have two different models designed using two set of dialogs: the *small model* is designed using 5 training conversations and the *medium model* has been designed adding other 10 conversation examples to the small. From the diagram model we obtained two sets of rules: $bs\_rules$ for the production of the belief state $B_t$ and $action\_rules$ for the production of the system action $A_t$. We also used $bs\_rules$ in two different configurations, that is, with or without the use of constraint on the previous action $A_{t-1}$ and we used $action\_rules$ in two different configurations, that is, with or without the constraint on the belief $B_t$.

We evaluated CLINN and the DAMD architecture [14] to determine their ability to recreate the inner states: the action span $A_t$ and the belief span $B_t$ as we aim to verify that our model can control the flow in the dialog states. To evaluate the ability to replicate $A_t$, we used the F1-measure that is the harmonic mean of recall and precision of produced actions with respect to gold actions. For what concerns the belief span we used the Joint Goal Accuracy that is the percentage of turns in a dialogue where the

---

[5] For an exhaustive description of the dialogue modeling please refer to [3]

user's informed joint goals are identified correctly. Joint goals are accumulated turn goals up to the current dialog turn.

| | | Injection Type | | | | | Action Span | Belief Span |
|---|---|---|---|---|---|---|---|---|
| System | Rule Set | Belief | Action | Action/Belief | Train Set | Test Set | F1 | joint goal |
| DAMD | | | | gold | 150 | full | 36.5 | 69.4 |
| CLINN | small | | no belief | gold | 150 | full | 39.8 | **71.9** |
| CLINN | small | | use belief | gold | 150 | full | 39.5 | 62.6 |
| CLINN | small | no action | | gold | 150 | full | 37.9 | 66.2 |
| CLINN | small | use action | | gold | 150 | full | **44.1** | 66.9 |
| DAMD | | | | gold | 1200 | full | 42.2 | 75.9 |
| CLINN | small | | no belief | gold | 1200 | full | 37.2 | 78.1 |
| CLINN | medium | | no belief | gold | 1200 | full | **47.2** | **82.4** |
| DAMD | | | | gen | 150 | full | 37.3 | 40.6 |
| CLINN | small | | no belief | gen | 150 | full | 39.6 | 54.3 |
| CLINN | small | | use belief | gen | 150 | full | 39.4 | 42.1 |
| CLINN | small | no action | | gen | 150 | full | 37.7 | 48.6 |
| CLINN | small | use action | | gen | 150 | full | **45.3** | **48.9** |
| DAMD | | | | gen | 1200 | full | 42.9 | 64 |
| CLINN | small | | no belief | gen | 1200 | full | 36.8 | 64.7 |
| CLINN | medium | | no belief | gen | 1200 | full | **48.8** | **69.4** |
| DAMD | | | | gen | 150 | reduced | 44.4 | 71.8 |
| DAMD | | | | gen | 1200 | reduced | 41.4 | 71.1 |
| CLINN | medium | | no belief | gen | 150 | reduced | 48.7 | 74.6 |
| CLINN | medium | | no belief | gen | 1200 | reduced | **53.4** | **84.5** |

**Table 1.** Comparison of the performances of DAMD and the CLINN system with different configurations. The type gold or gen in Action/Belief denotes if previous Action/Belief are taken from the ground truth (gold) or are generated by the system (gen).

### 3.2 Results and discussion

The first set of the experimental results (Table 1 - Test Set "Full") shows that CLINN positively inject symbolic rules in sequence-to-sequence neural networks when training data are scarce. CLINN outperforms DAMD in nearly all the configurations when compared on the Action Span F1 and in some configuration when compared on the joint goal on the Belief Span. More importantly, CLINN seems to obtain interesting results in situations with data scarcity. With a small training set with 150 dialogs, one configuration of CLINN outperforms DAMD of more than 7.5% on the Action Span F1 both in the *gold* setting (44.1 vs. 36.5) and in the *gen* setting (45.3 vs. 37.5). The increase in the joint goal for the Belief Span is less impressive in the *gold* setting where only one configuration – with rule injection type Action without using belief constraints – outperforms DAMD (71.9 vs. 69.4). Instead, the performance increase of CLINN in the joint goal is more stable in the *gen* setting. Moreover, the difference between DAMD

and the best system is more than 13% (54.3 vs. 40.6). Moreover, CLINN is an effective model to include hand-crafted rules when the training set is relatively large. We selected the best configuration selected with the training set of 150 dialogs (Injection Type Action with no belief) and we experimented with 1,200 dialogs as training. By using a larger rule set, that is, the *medium* rule set, CLINN outperforms DAMD for the action spans and for the joint goal of the belief span in the *gold* and in *gen* setting.

The second set of experimental results (Table 1 - Test Set "reduced") gives the important indication that CLINN can help in controlling the behavior of dialog systems in specific and critical situations. The reduced test set is composed only with the conversations used for building the medium rule set (15 conversations). Although the DAMD model contains these conversations in the training set, its performance drops when increasing the training set. CLINN instead improves its performance of both metrics when the training set increases. Hence, CLINN offer a better stability for critical dialogs that are used to design rules.

The two sets of experiments demonstrates the applicability of CLINN on industrial real cases.

## 4   Conclusions

Critical industrial applications such as banking or medical applications impose important constraints on Conversational AI systems: data scarcity and need for certified dialogs. We proposed Conversational Logic Injected Neural Network that allow to positively include logical rules to control a sequence-to-sequence dialog manager. Our system shows a possible approach towards a more effective integration of neural network conversational AI in industrial applications.

# References

1. Budzianowski, P., Wen, T.H., Tseng, B.H., Casanueva, I., Ultes, S., Ramadan, O., Gašić, M.: MultiWoz - A large-scale multi-domain wizard-of-oz dataset for task-oriented dialogue modelling. In: Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, EMNLP 2018 (2020)

2. Dinan, E., Roller, S., Shuster, K., Fan, A., Auli, M., Weston, J.: Widzard Of Wikipedia: Knowledge-powered conversational agents. In: 7th International Conference on Learning Representations, ICLR 2019 (2019)

3. Giannone, C., Bellomaria, V., Favalli, A., Romagnoli, R.: Iride R : an Industrial Perspective on Production Grade End To End Dialog System. In: Proceeting of the Italian Conference of Computational Linguistics (CLIC). Bari (2019), `https://www.gartner.com/smarterwithgartner/4-trends-`

4. Henderson, M., Vulic, I., Gerz, D., Casanueva, I., Budzianowski, P., Coope, S., Spithourakis, G., Wen, T.H., Mrkšic, N., Su, P.H.: Training neural response selection for task-oriented dialogue systems. In: ACL 2019 - 57th Annual Meeting of the Association for Computational Linguistics, Proceedings of the Conference. pp. 5392–5404 (6 2019), `http://arxiv.org/abs/1906.01543`

5. Lei, W., Jin, X., Ren, Z., He, X., Kan, M.Y., Yin, D.: Sequicity: Simplifying task-oriented dialogue systems with single sequence-to-sequence architectures. In: ACL 2018 - 56th Annual Meeting of the Association for Computational Linguistics, Proceedings of the Conference (Long Papers). vol. 1, pp. 1437–1447. Association for Computational Linguistics (2018). https://doi.org/10.18653/v1/p18-1133, `http://github.com/WING-NUS/sequicity`

6. Lison, P., Kennington, C.: OpenDial: A toolkit for developing spoken dialogue systems with probabilistic rules. In: 54th Annual Meeting of the Association for Computational Linguistics, ACL 2016 - System Demonstrations. pp. 67–72 (2016). https://doi.org/10.18653/v1/p16-4012, `http://www.opendial-toolkit.net`

7. Serban, I.V., Sordoni, A., Bengio, Y., Courville, A., Pineau, J.: Building end-To-end dialogue systems using generative hierarchical neural network models. 30th AAAI Conference on Artificial Intelligence, AAAI 2016 pp. 3776–3783 (2016)

8. Shah, P., Hakkani-Tür, D., Tür, G., Rastogi, A., Bapna, A., Nayak, N., Heck, L.: Building a Conversational Agent Overnight with Dialogue Self-Play (1 2018), `http://arxiv.org/abs/1801.04871`

9. Sordoni, A., Bengio, Y., Vahabi, H., Lioma, C., Simonsen, J.G., Nie, J.Y.: A hierarchical recurrent encoder-decoder for generative context-aware query suggestion. In: International Conference on Information and Knowledge Management, Proceedings. vol. 19-23-Oct-, pp. 553–562. Association for Computing Machinery, New York, New York, USA (10 2015). https://doi.org/10.1145/2806416.2806493, `http://dl.acm.org/citation.cfm?doid=2806416.2806493`

10. Vlasov, V., Mosig, J.E.M., Nichol, A.: Dialogue Transformers (2019), `http://arxiv.org/abs/1910.00486`

11. Wen, T.H., Su, P.H., Budzianowski Pawełand Casanueva, I., Vulić, I.: Data Collection and End-to-End Learning for Conversational AI. Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP): Tutorial Abstracts (2019)

12. Williams, J.D., Asadi, K., Zweig, G.: Hybrid code networks: Practical and efficient end-to-end dialog control with supervised and reinforcement learning. In: ACL 2017 - 55th Annual Meeting of the Association for Computational Linguistics, Proceedings of the Conference (Long Papers). vol. 1, pp. 665–677 (2017). https://doi.org/10.18653/v1/P17-1062

13. Zanzotto, F.M.: Viewpoint: Human-in-the-loop Artificial Intelligence. J. Artif. Intell. Res. **64**, 243–252 (2019). https://doi.org/10.1613/jair.1.11345, `https://doi.org/10.1613/jair.1.11345`

14. Zhang, Y., Ou, Z., Yu, Z.: Task-Oriented Dialog Systems that Consider Multiple Appropriate Responses under the Same Context. Proceedings of the AAAI Conference on Artificial Intelligence **34**(05), 9604–9611 (4 2019). https://doi.org/10.1609/aaai.v34i05.6507, `www.aaai.orghttp://arxiv.org/abs/1911.10484https://www.aaai.org/ojs/index.php/AAAI/article/view/6507`