

---

# Deep Reinforcement Learning for Multi-Domain Dialogue Systems\*

---

Heriberto Cuayáhuitl<sup>1</sup>, Seunghak Yu<sup>2</sup>, Ashley Williamson<sup>1</sup>, Jacob Carse<sup>1</sup>

<sup>1</sup>University of Lincoln, School of Computer Science, United Kingdom

<sup>2</sup>Artificial Intelligence Team, Samsung Electronics Co. Ltd., Seoul, South Korea  
HCuayahuitl@lincoln.ac.uk

## Abstract

Standard deep reinforcement learning methods such as Deep Q-Networks (DQN) for multiple tasks (domains) face scalability problems. We propose a method for multi-domain dialogue policy learning—termed NDQN, and apply it to an information-seeking spoken dialogue system in the domains of restaurants and hotels. Experimental results comparing DQN (baseline) versus NDQN (proposed) using simulations report that our proposed method exhibits better scalability and is promising for optimising the behaviour of multi-domain dialogue systems.

## 1 Introduction

Dialogue systems based on the Reinforcement Learning (RL) paradigm offer the possibility to treat dialogue design as an optimisation problem, and are attractive because they can improve their performance over time with experience. But the application of RL is not trivial due to the complexity of the problem such as large state-action spaces exhibited in human-machine conversations. This is especially true in multi-domain systems, where the number of state variables (features) and dialogue actions increases rapidly as more domains are taken into account. On the one hand, unique situations in the interaction can be described by a large number of variables (e.g. words raised in the conversation by the system and user) so that enumerating them would result in very large state spaces. On the other hand, the action space can also be large due to the wide range of unique dialogue actions (e.g. requests, apologies, confirmations in multiple contexts).

While one can aim to optimise the interaction via compression of the search space, it is usually unclear what features to incorporate in the state representation. This is a strong motivation for applying Deep Reinforcement Learning (DRL) to dialogue management so that the agent can simultaneously learn its feature representation and policy [15]. This paper makes use of raw noisy text as features in an attempt to avoid engineered features to represent the dialogue state. By using this representation, dialogue agents bypass spoken language understanding in order to learn dialogue policies directly from raw (noisy) text to actions [1].

We address dialogue optimisation using the divide-and-conquer approach, in which dialogue states can be described at different levels of granularity, and an action can execute behaviour using either a single dialogue action (taking one dialogue turn) or a composite one (equivalent to a subdialogue taking multiple dialogue turns). This approach offers at least two benefits: (a) modularity helps to optimise subdialogues that may be easier to optimise than the whole dialogue; and (b) subdialogues may include only relevant dialogue knowledge in the states and relevant actions, thus reducing significantly the size of possible solutions: consequently they can be found faster. These properties

---

\*Funding from Samsung Electronics Ltd. and the University of Lincoln is gratefully acknowledged. We thank Raymond Kirk for helping with App development (to integrate the agents in this paper) and system testing. We also thank Heesik Jeon and the AI team at Samsung for their executive efforts in this project.

are crucial for training the behaviour of multi-domain spoken dialogue systems in which there may be a large set of state variables or a large number of actions.

Below we describe a data-driven method to the approach described, which we have applied to an information-seeking dialogue system in the domains of restaurants and hotels. Experimental results show that the proposed method can train policies faster and more effectively than a standard algorithm in the literature, showing promise for training multi-domain dialogue systems.

## 2 Literature Review

Recently, multi-domain spoken conversational agents have received an increasing amount of attention. This may be due to the fact that speech technologies such as Automated Speech Recognition (ASR) and Text-To-Speech (TTS) have reached a high degree of maturity. But the question of *How to design conversational systems for human-machine interaction in multiple domains (or tasks)?* is still an open and interesting problem in artificial intelligence. The dialogue system proposed by [11] used a distributed architecture of domain experts modulated by a domain selector. The latter used a decision tree with classification errors over 20% in 5 domains. This indicates that not only individual domains have to exhibit robust interactions against errors, but also that errors increase by incorporating more domains.[9] used rule-based classifiers for predicting user intentions, which are executed using a Hierarchical Task Network (HTN) incorporating expert knowledge. Trainable multi-domain dialogue systems using traditional reinforcement learning include [4, 13, 22, 5]. These systems use a modest amount of features, and in contrast to neural-based systems, they require manual feature engineering.

Recent work on deep learning applied to task-oriented conversational agents include the following. [6] uses a Recurrent Neural Network (RNN) for dialogue act prediction in a POMDP-based dialogue system, which focuses on mapping system and user sentences to dialogue acts. [2] applies Deep Reinforcement Learning with a fully-connected neural network for trading negotiations in board games, which focuses on mapping game situations to dialogue actions. [20] trains RNN-based classifiers for predicting dialogue success in multi-domain dialogue systems, which can be applied to unseen domains. [17] also trains RNN-based classifiers but for belief tracking in order to improve the robustness of recognised user responses across dialogue turns. Other neural-based conversational agents have been applied to text prediction using the sequence-to-sequence approach [19, 21], and to reasoning with inference for text-based question answering [23].

From these works, we can observe that supervised learning is the dominating form of training in neural-based conversational agents. To our knowledge, we report the first multi-domain dialogue system using deep reinforcement learning. This form of learning is interesting because it can perform feature learning and policy learning simultaneously, and its effective application in real-world dialogue scenarios remains to be demonstrated.

## 3 Method

Our proposed method to scale up Deep Reinforcement Learning (DRL) for multi-domain dialogue systems has two stages. First, multi-policy learning via a network of DRL agents; and second, more compact state representations by compressing raw inputs. Although these two stages can be applied independently, their combination aims for further scalability than any one of them individually.

### 3.1 Network of Deep Q-Networks

We propose to optimise multi-domain dialogue systems using a network of Deep Reinforcement Learners, for example a network of Deep Q-networks (DQN) — see [15, 16] for an introduction to the standard DQN method. In our method, instead of training a single DQN, we train a set of DQNs (also referred to as NDQNs), where every DQN represents a specialised skill to converse in a particular subdialogue — see Figure 1. In addition, the network of agents enable DQNs to be executed without a fixed structure in order to support flexible and unstructured dialogues. In contrast to hierarchical DQNs [12] that follow a strict sequence of agents, in our method an NDQN allows transitions between all DQN agents except for self-transitions and loops (the latter using a stack-based approach as in [3]). Furthermore, while user responses can motivate transitions to

another domain in the network, completing a subdialogue within a domain motivates a transition to the previous domain to resume the interaction. Algorithm 1 describes the procedure to train and execute NDQNs.

An optimal policy in an NDQN performs action selection according to

$$\pi_{\theta^{(d)}}^*(s) = \arg \max_{a \in A^{(d)}} Q^{*(d)}(s, a; \theta^{(d)}), \quad (1)$$

where domain or skill  $d \in D$  is selected according to

$$d = \arg \max_{d' \in D} F(d'|d, \mathbf{e}), \quad (2)$$

and evidence  $\mathbf{e}$  takes into account all features that describe the environment space of domain  $d$ . While this transition function (Eq. 2) is used for high-level transitions in the interaction, Equation 1 is used for low-level transitions within a node (skill) in the network and subject to reinforcement learning. NDQN assumes that the domain transition function  $F$  can be deterministic or probabilistic (the latter due to uncertainty in the interaction), and it is a prior requirement for NDQN-Learning.

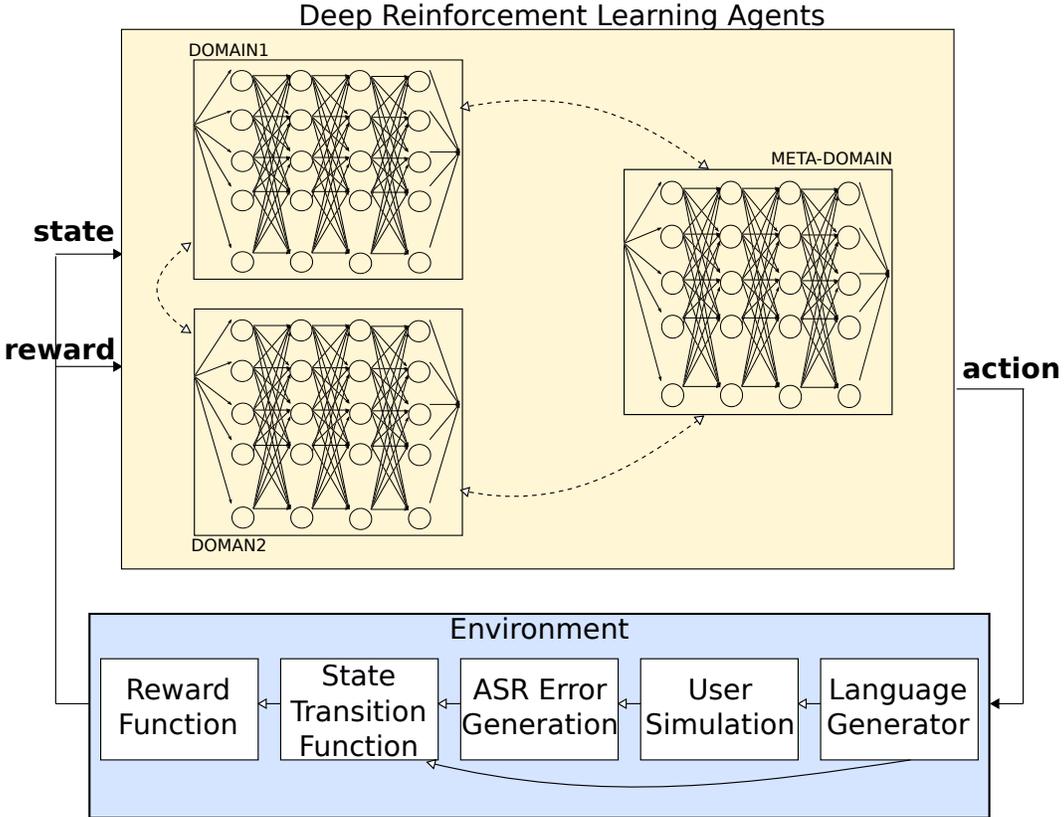


Figure 1: Multi-domain DRL agents with flexible interaction. The dashed arrows connecting domains denote flexible transitions between domains in order to avoid a rigid structure in the interaction. Although all policies are considered for decision-making, only one domain can be executed at a time implying that a previous domain can continue its execution in order to resume the interaction.

### 3.2 NDQNs with Compressed Raw Inputs

Previous work on dialogue policy learning using DRL map raw (noisy) text to actions [1, 24]. This is not only computationally intensive, but it becomes infeasible for dialogue systems with large vocabularies. To tackle this problem we propose to use delexicalised sentences (as proposed in [8]) and synonymised sentences. This has the advantage that dialogue policies can be trained from more compact state representations than those using only raw inputs, and have coverage for a larger vocabulary than trained for.

---

**Algorithm 1** Network of Deep Q-Learners (NDQN)

---

- 1: Initialise set of Deep Q-Networks with replay memories  $D^{(d)}$ , action-value functions  $Q^{(d)}$  with random weights  $\theta^{(d)}$ , and target action-value functions  $\hat{Q}^{(d)}$  with weights  $\hat{\theta}^{(d)} = \theta^{(d)}$
  - 2: **repeat**
  - 3:    $d \leftarrow$  initial domain, predefined or defined by  $\arg \max_{d \in D} F_o(d)$
  - 4:    $s \leftarrow$  initial environment state in  $S^{(d)}$
  - 5:   **repeat**
  - 6:     **repeat**
  - 7:       Choose action  $a \in A^{(d)}$  in  $s$  derived from  $Q^{(d)}$  (e.g.  $\epsilon$ -greedy, Thompson)
  - 8:       Execute action  $a$  and observe reward  $r$  and next state  $s'$
  - 9:       Append transition  $(s, a, r, s')$  to  $D^{(d)}$
  - 10:       $B^{(d)} \leftarrow$  sample random minibatch of experiences from  $D^{(d)}$
  - 11:       $d' \leftarrow$  select next domain according to  $\arg \max_{d' \in D} F(d'|s', \mathbf{e})$
  - 12:      
$$y_j = \begin{cases} r_j & \text{if final step of episode} \\ r_j + \gamma \max_{a \in A^{(d)}} \hat{Q}^{(d)}(s', a'; \hat{\theta}^{(d)}), & \text{otherwise} \end{cases}$$
  - 13:      Gradient descent step on  $(y_j - Q^{(d)}(s', a'; \theta^{(d)}))^2$  using  $B^{(d)}$
  - 14:      Reset  $\hat{Q}^{(d)} = Q^{(d)}$  every  $C$  steps
  - 15:       $s \leftarrow s'$
  - 16:     **until**  $s$  is a terminal state or  $d \neq d'$
  - 17:      $d \leftarrow d'$
  - 18:   **until**  $s$  is a goal state
  - 19: **until** convergence
- 

### 3.2.1 Delexicalisation

Consider a dialogue system for restaurant search receiving the following user request—with corresponding delexicalised sentence underneath.

I am looking for **italian** food in the **city centre**  
↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓  
I am looking for **\$foodtype** food in the **\$area**

The latter representation combining words and slot IDs (denoted with the symbol ‘\$’) has several practical advantages. For example, policies can be learnt faster, they contribute to further scalability of systems with large vocabularies, and policies do not have to be retrained if the slot values change over time. In this work we use heuristics to replace slot values by slot IDs, and a trainable component for automatic slot labelling is considered beyond the scope of this paper.

### 3.2.2 Synonymization

Consider the same system above receiving the following user request given the unknown words ‘fancy’ and ‘cuisine’—with corresponding synonyms underneath.

We **fancy** italian **cuisine** in the centre of town  
↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓  
We **want** italian **food** in the centre of town

We argue that word synonyms can be useful in such situations because the unknown word ‘fancy’ can trigger the known word feature ‘want’. Similarly, the unknown word ‘cuisine’ can trigger the known word feature ‘food’. In this way, the vocabulary of our NDQNs incorporate a mapping from filler words and slot values to synonyms in order to cope with unseen wordings. Due to the complexity of automatic generation of meaningful synonyms, our synonyms have been manually specified but they can be generated automatically for example from word embeddings [14].

## 4 Multi-Domain Dialogue System

The proposed computational framework for training multi-domain dialogue agents is a substantial extension from the publicly available software tools SimpleDS [1] and ConvnetJS [10]. It can be executed in training or test mode using simulations or speech-based interactions (via an App). Our dialogue system runs under a client-server architecture, where the learning agents—one per domain—act as the *clients* and the dialogue system as the *server*. They communicate by exchanging messages, where the clients tell the server the action to execute, and the server tells the clients the state and reward observed. The elements for training multi-domain dialogue systems are as follows.

**State Spaces** They include word-based features depending on the vocabulary of each learning agent. They include 177 unique words<sup>2</sup> without synonyms, and 150 unique words with synonyms. For example, an agent in the domain of restaurants has relevant features for its domain and it is agnostic of features in other domains. While words derived from system responses are treated as binary variables (i.e. word present or absent), the words derived from noisy user responses can be seen as continuous variables by taking ASR confidence scores into account. Since a single variable per word is used, user features override system ones in case of overlaps.

**Action Spaces** They include dialogue acts for the targeted domains—currently 69 unique actions in total. Example dialogue act types, dialogue acts without parameters, are as follows: Salutation(), Request(), AskFor(), Apology(), ExpConfirm(), ImpConfirm(), Retrieve(), Provide(), among others. Rather than learning with whole action sets, our framework supports learning from constrained actions by applying learning updates only on the set of valid actions. These actions are derived from the most likely actions,  $Pr(a|s) > 0.0001$ , from Naive Bayes classifiers (due to scalability purposes) trained from example dialogues. See example demonstration dialogue in Appendix A. In addition to the most probable data-like actions, the constrained actions are extended with the legitimate requests, apologies and confirmations. The fact that constrained actions are data-driven and driven by domain-independent heuristics, facilitates its usage across domains.

**State Transition Functions** They are based on numerical vectors representing the last system and user responses. Taking a wider dialogue context is also possible but not explored in this paper. The system responses are straightforward, 0 if absent and 1 if present (hit-or-miss). The user responses correspond to the confidence level [0..1] of noisy user responses. While system responses are generated from stochastic templates, user responses are generated from semi-random user behaviour. These elements enable the creation of a vast amount of different conversations for agent training.

**Domain Transition Function** This function specifies the next domain or task in focus. It is currently defined deterministically, and it is also implemented as a SVM classifier trained from example interactions—see Appendix A. The design of this classifier follows that of a two-deep fully connected neural network with 80 nodes in each hidden layer, with tanh activation, and an SVM output layer, using Hinge Loss. While the input layer accepts domain-independent *words-as-features* vectors representing the unique global vocabulary shared amongst all domains in a hit-or-miss approach, the output layer has 3 classes representing system domains (meta<sup>3</sup>, restaurants and hotels). 15K dialogues of data were generated, partitioned as a 60-40 training-testing split, and trained for 180 epochs. Initial results of this classifier shows a 87.5% classification accuracy on user-simulated data.

**Reward Function** It is defined as  $R(s, a, s') = GR + DR - DL$ , where  $GR$  is goal-based reward treated as task success [0..1] (the proportion of positively confirmed slots and information retrieved and presented);  $DR$  is a data-like probability of having observed action  $a$  in state  $s$ ; and  $DL = t * w$  is a dialogue length measure used to encourage efficient interactions with  $t$  time steps and weight  $w$  (-0.1 in our case). The  $DR$  scores are derived from Naive Bayes classifiers to allow statistical inference over actions given states ( $Pr(a|s)$ ).

---

<sup>2</sup>The unique words in our system’s vocabulary excludes words from information presentation due to the vast amount of information about hotels and restaurants. Nonetheless and during testing, our system retrieves live information from <http://www.bookatable.co.uk> and [www.reservetravel.com](http://www.reservetravel.com).

<sup>3</sup>We refer to meta domain as subdialogues containing domain-general system and user responses.

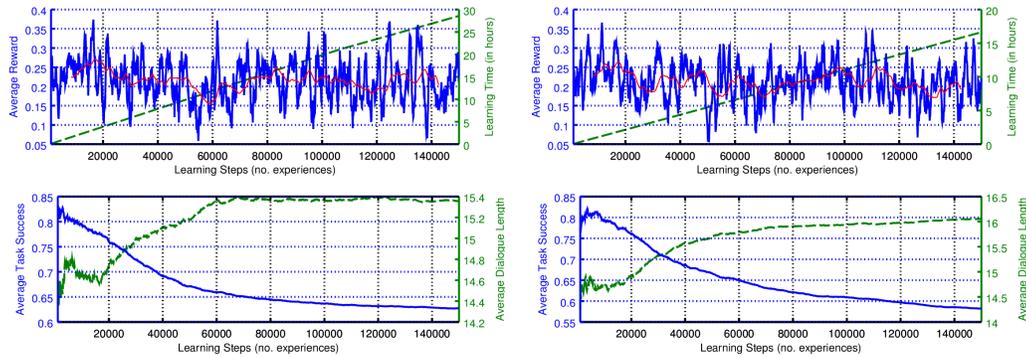


Figure 2: Learning curves of the baseline system: (left) without input compression, (right) with input compression. The higher the better in blue straight lines, and the lower the better in other metrics.

**Model Architectures** We use fully-connected multilayer neural nets, trained with stochastic gradient descent, where nodes in the input layers depend on the vocabulary of each agent. The use of convolutional neural nets is work in progress. They include 2 hidden layers with 80 nodes with Rectified Linear Units to normalise their weights [18]. Dropout[7] and adaptive learning rates are also part of our work in progress. Other hyperparameters include experience replay size=10000, burning steps=1000, discount factor=0.7, minimum epsilon=0.001, batch size=32, and learning steps=30000.

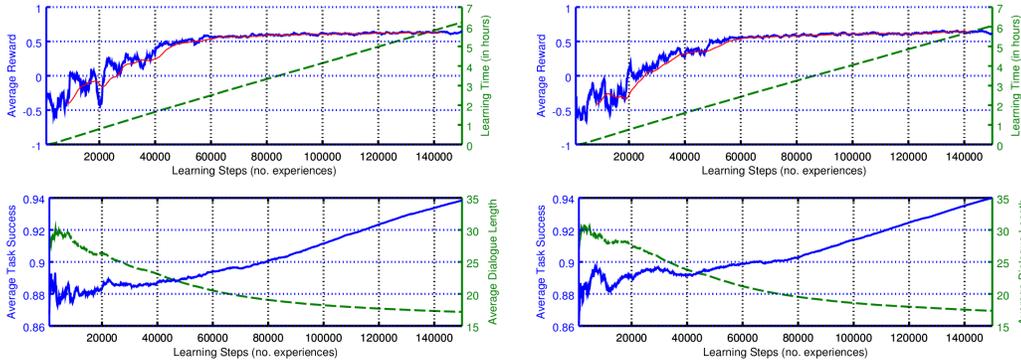
## 5 Experimental Results

Here we compare a multi-domain dialogue system using a standard DRL method versus our proposed method described in Sections 3 and 4. While the former (DQN) uses a single policy for learning (*baseline*), the latter (NDQN) uses multiple policies (*proposed*). Both multi-domain dialogue systems use the same data, resources and hyperparameters for training, the only difference between both systems is the learning method (DQN or NDQN) or state representation (with or without compression).

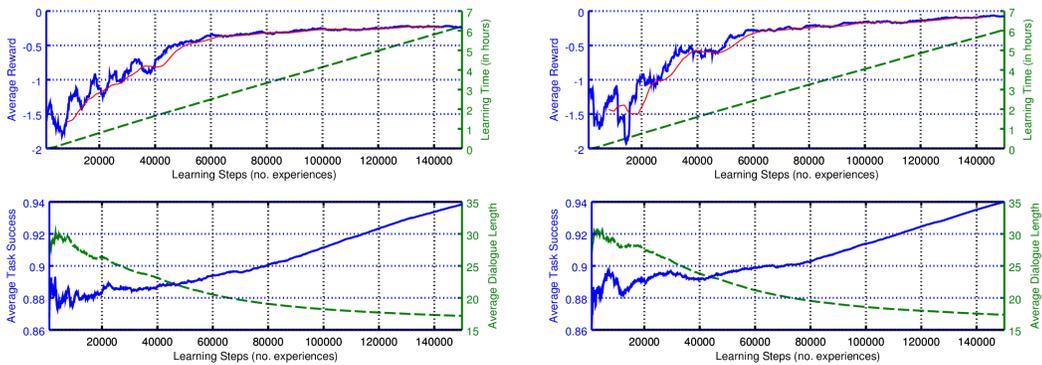
We use four different metrics to measure system performance: avg. reward, learning time, avg. task success, and avg. dialogue length. The higher the better in the first and third, and the lower the better in the second and fourth. Figure 2 shows learning curves for the baseline DQN-based system, and Figure 3 shows learning curves for the proposed NDQN-based system. Both the baseline and proposed system report results over 150K learning steps (about 8700 dialogues). Our results report that training multi-domain systems using a single policy is twofold harder than using a multi-policy approach. First, this is evidenced by the fact that the baseline policies do not improve over time, and the policies with the proposed method do. This is presumably due to the abstraction exhibited in the multi-policy approach—more focused system actions rather than interleaving them across domains. Second, our proposed system also learned 4.6 times faster than the baseline, which was accelerated further to 4.7 times faster by using compressed inputs<sup>4</sup>. By applying synonymization we are able to use a smaller vocabulary when training and then a much wider vocabulary at runtime, which adds robustness in the presence of unseen dialogues. These results show indication of better scalability for NDQN to multiple domains.

Although the currently generated dialogues using the proposed method seem reasonable, a natural question to ask is *How good (qualitatively speaking) are the trained policies?* This question will be answered in an evaluation reported in future work.

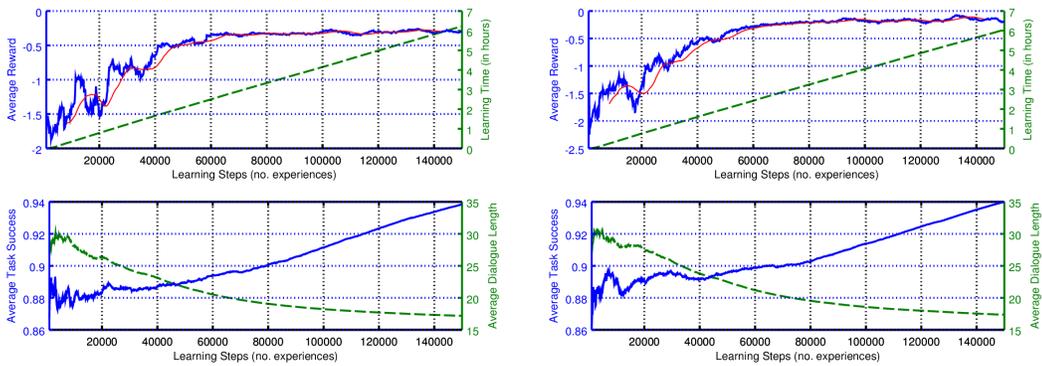
<sup>4</sup>Ran on Intel Core i5-3210M CPU @ 2.50GHz x 4; 8GiB DDR4 RAM @ 2400MHz.



(a) Meta Domain: (left) without input compression, (right) with input compression



(b) Restaurants Domain: (left) without input compression, (right) with input compression



(c) Hotels Domain: (left) without input compression, (right) with input compression

Figure 3: Learning curves of the proposed system using simultaneous policy learning with the proposed method. The higher the better in avg. reward and avg. task success, and the lower the better in other metrics. The plots on the left correspond to our proposed system with word-based features, and the plots on the right correspond to our proposed system with delexicalised inputs as features. The latter plots show no performance degradation despite of using more compact state representations.

Method / System	Baseline (DQN)	Proposed (NDQN)
Without input compression	28.57 hrs	6.21 hrs
With input compression	16.63 hrs	6.05 hrs

Table 1: Learning times of the baseline and proposed systems

## 6 Conclusion and Future Work

We have described a method for training multi-domain dialogue systems in a more scalable way than traditional deep reinforcement learning, e.g. using the DQN method. The proposed method uses a Network of DQN (NDQN) agents in order to train specialised agents, and compression of input features. Experimental results using simulations report that the proposed method (NDQN) can train policies faster and more effectively than DQN—almost 5 times faster than DQNs. In addition, our results report that input compression contribute faster learning without performance degradation across metrics. These dialogue agents are the first applications of deep reinforcement learning to optimising the behaviour of multi-domain human-machine dialogues.

In future work we plan to carry out a data collection of speech-based interactions with recruited participants in order to evaluate and/or retrain our multi-domain dialogue policies with a larger dataset. Other future work includes applying and refining the proposed framework to larger sets of domains, and applying it to multimodal interaction in different types of devices and machines.

### A Example Demonstration Dialogue

We aim for a data-efficient system, trainable from a small set of demonstrations. This dialogue is one of the 12 manually-specified seed demonstrations for inducing the system’s features (words) and actions (dialogue acts), and for training classifiers to predict the most likely actions given the previous system and user responses. The latter were used by the DRL agents described in Section 3 for inducing their dialogue behaviour according to the learning environment described in Section 4. Square brackets denote system actions, double quotes denote system or user responses, and words beginning with an underscore encapsulate live web-based information according to the user requests.

Domain	Agent	[Action]"Verbalisation"
Meta	SYS	[Salutation(greeting)]"Welcome to MultiDS!"
	USR	" <i>I am looking for a hotel in edinburgh on the 2nd of january for 2 nights</i> "
Hotels	SYS	[Apology(h_nights)]"Sorry, how long?"
	USR	" <i>2 nights</i> "
	SYS	[ExpConfirm(h_day=\$h_day,h_month=\$h_month,h_nights=\$h_nights)]
	USR	" <i>yes i did</i> "
	SYS	[ExpConfirm(h_city=\$h_city)]"Did you say in edinburgh?"
	USR	" <i>i did</i> "
	SYS	[Retrieve(h_info)]"Let me see." [Provide(h_info)]"I found the following hotels. <u>hotels</u> ."
Meta	USR	[AskFor(h_more)]"Anything else?" " <i>I am looking for cheap japanese food in the centre of town</i> "
Restaurants	SYS	[Apology(r_area)]"Sorry, what area?"
	USR	" <i>In the centre</i> "
	SYS	[ExpConfirm(r_price=\$r_price,r_area=\$r_area)]
	USR	" <i>yes i did</i> "
	SYS	[ExpConfirm(r_food=\$r_food)]"Did you say japanese food?"
	USR	" <i>i did</i> "
Meta	SYS	[Retrieve(h_info)]"Let me see." [Provide(h_info)]"I found the following restaurants. <u>res</u> ."
	USR	Anything else?" " <i>no</i> "
	SYS	[Salutation(closing)]"Okay, talk to you soon. Bye!"

## References

- [1] H. Cuayáhuitl. SimpleDS: A simple deep reinforcement learning dialogue system. *CoRR*, abs/1601.04574, 2016.
- [2] H. Cuayáhuitl, S. Keizer, and O. Lemon. Strategic dialogue management via deep reinforcement learning. *CoRR*, abs/1511.08099, 2015.
- [3] H. Cuayáhuitl, I. Kruijff-Korbayová, and N. Dethlefs. Nonstrict hierarchical reinforcement learning for interactive systems and robots. *ACM Transactions on Interactive Intelligent Systems (TiiS)*, 4(3), 2014.
- [4] H. Cuayáhuitl, S. Renals, O. Lemon, and H. Shimodaira. Evaluation of a hierarchical reinforcement learning spoken dialogue system. *Computer Speech & Language*, 24(2), 2010.
- [5] M. Gasic, N. Mrksic, P. Su, D. Vandyke, T. Wen, and S. J. Young. Policy committee for adaptation in multi-domain spoken dialogue systems. In *ASRU*, 2015.
- [6] W. Ge and B. Xu. Dialogue management based on multi-domain corpus. In *Annual Meeting of the Special Interest Group on Discourse and Dialogue (SIGDIAL)*, 2015.
- [7] I. J. Goodfellow, D. Warde-Farley, M. Mirza, A. C. Courville, and Y. Bengio. Maxout networks. In *International Conference on Machine Learning (ICML)*, 2013.
- [8] M. Henderson, B. Thomson, and S. J. Young. Robust dialog state tracking using delexicalised recurrent neural networks and unsupervised adaptation. In *IEEE Spoken Language Technology Workshop*, 2014.
- [9] H. Jeon, H. R. Oh, I. Hwang, and J. Kim. An intelligent dialogue agent for the IoT home. In *AAAI Workshop on AI Applied to Assistive Technologies and Smart Environments*, 2016.
- [10] A. Karpathy. ConvNetJS: Javascript library for deep learning. <http://cs.stanford.edu/people/karpathy/convnetjs/>, 2015.
- [11] K. Komatani, N. Kanda, M. Nakano, K. Nakadai, H. Tsujino, T. Ogata, and H. G. Okuno. Multi-domain spoken dialogue system with extensibility and robustness against speech recognition errors. In *SIGDial Workshop on Discourse and Dialogue*, 2006.
- [12] T. D. Kulkarni, K. Narasimhan, A. Saeedi, and J. B. Tenenbaum. Hierarchical deep reinforcement learning: Integrating temporal abstraction and intrinsic motivation. *CoRR*, abs/1604.06057, 2016.
- [13] P. Lison. Multi-policy dialogue management. In *SIGDIAL*, 2011.
- [14] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, and J. Dean. Distributed representations of words and phrases and their compositionality. In *NIPS*, 2013.
- [15] V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. Riedmiller. Playing atari with deep reinforcement learning. In *NIPS Deep Learning Workshop*. 2013.
- [16] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski, S. Petersen, C. Beattie, A. Sadik, I. Antonoglou, H. King, D. Kumaran, D. Wierstra, S. Legg, and D. Hassabis. Human-level control through deep reinforcement learning. *Nature*, 518(7540), 02 2015.
- [17] N. Mrksic, D. Ó. Séaghdha, B. Thomson, M. Gasic, P. Su, D. Vandyke, T. Wen, and S. J. Young. Multi-domain dialog state tracking using recurrent neural networks. *CoRR*, abs/1506.07190, 2015.
- [18] V. Nair and G. E. Hinton. Rectified linear units improve restricted boltzmann machines. In *International Conference on Machine Learning (ICML)*, 2010.
- [19] I. V. Serban, A. Sordoni, Y. Bengio, A. C. Courville, and J. Pineau. Hierarchical neural network generative models for movie dialogues. *CoRR*, abs/1507.04808, 2015.
- [20] D. Vandyke, P. Su, M. Gasic, N. Mrksic, T. Wen, and S. J. Young. Multi-domain dialogue success classifiers for policy training. In *ASRU*, 2015.
- [21] O. Vinyals and Q. V. Le. A neural conversational model. *CoRR*, abs/1506.05869, 2015.
- [22] Z. Wang, H. Chen, G. Wang, H. Tian, H. Wu, and H. Wang. Policy learning for domain selection in an extensible multi-domain spoken dialogue system. In *Empirical Methods in Natural Language Processing EMNLP*, 2014.
- [23] J. Weston, S. Chopra, and A. Bordes. Memory networks. *CoRR*, abs/1410.3916, 2014.
- [24] T. Zhao and M. Eskénazi. Towards end-to-end learning for dialog state tracking and management using deep reinforcement learning. *CoRR*, abs/1606.02560, 2016.