# Deep Reinforcement Learning for Chatbots Using Clustered Actions and Human-Likeness Rewards

Heriberto Cuayáhuitl
*School of Computer Science*
*University of Lincoln*
Lincoln, United Kingdom
HCuayahuitl@lincoln.ac.uk

Donghyeon Lee
*Artificial Intelligence Research Group*
*Samsung Electronics*
Seoul, South Korea
dh.semko.lee@samsung.com

Seonghan Ryu
*Artificial Intelligence Research Group*
*Samsung Electronics*
Seoul, South Korea
seonghan.ryu@samsung.com

Sungja Choi
*Artificial Intelligence Research Group*
*Samsung Electronics*
Seoul, South Korea
sungja.choi@samsung.com

Inchul Hwang
*Artificial Intelligence Research Group*
*Samsung Electronics*
Seoul, South Korea
inc.hwang@samsung.com

Jihie Kim
*Artificial Intelligence Research Group*
*Samsung Electronics*
Seoul, South Korea
jihie.kim@samsung.com

*Abstract*—Training chatbots using the reinforcement learning paradigm is challenging due to high-dimensional states, infinite action spaces and the difficulty in specifying the reward function. We address such problems using clustered actions instead of infinite actions, and a simple but promising reward function based on human-likeness scores derived from human-human dialogue data. We train Deep Reinforcement Learning (DRL) agents using chitchat data in raw text—without any manual annotations. Experimental results using different splits of training data report the following. First, that our agents learn reasonable policies in the environments they get familiarised with, but their performance drops substantially when they are exposed to a test set of unseen dialogues. Second, that the choice of sentence embedding size between 100 and 300 dimensions is not significantly different on test data. Third, that our proposed human-likeness rewards are reasonable for training chatbots as long as they use lengthy dialogue histories of $\geq 10$ sentences.

*Index Terms*—neural networks, reinforcement / unsupervised / supervised learning, sentence embeddings, chatbots, chitchat

## I. INTRODUCTION

What happens in the minds of humans during chatty interactions containing sentences that are not only coherent but also engaging? While not all chatty human dialogues are engaging, they are arguably coherent [1]. They also exhibit large vocabularies—according to the language in focus—because conversations can address any topic that comes to the minds of the partner conversants. In addition, each contribution by a partner conversant may exhibit multiple sentences instead of one such as greeting+question or acknowledgement+statement+question. Furthermore, the topics raised in the conversation may go back and forth without losing coherence. This is a big challenge for data-driven chatbots.

We present a novel approach based on the reinforcement learning [2], unsupervised learning [3] and deep learning [4] paradigms. Our learning scenario is as follows: given a data set of human-human dialogues in raw text (without any manually

Work carried out while the first author was visiting Samsung Research.
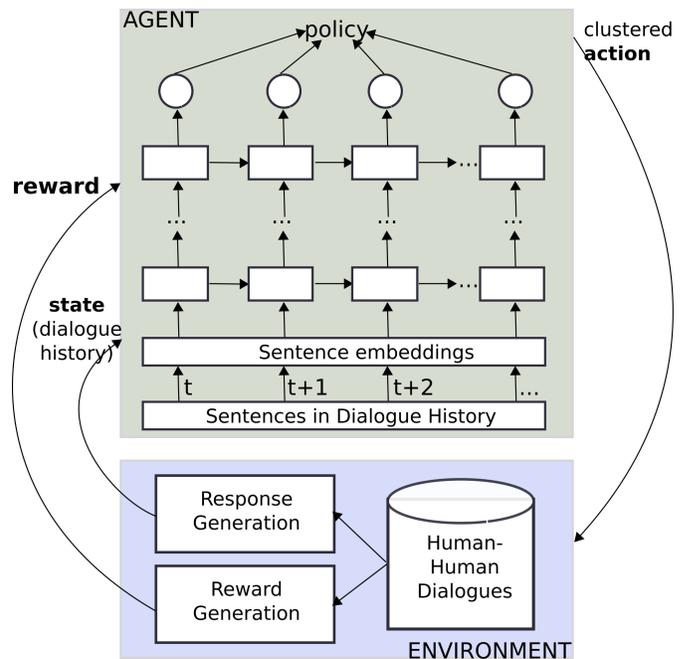


Fig. 1. High-level architecture of the proposed deep reinforcement learning approach for chatbots—see text for details

provided labels), a Deep Reinforcement Learning (DRL) agent takes the role of one of the two partner conversants in order to learn to select human-like sentences when exposed to both human-like and non-human-like sentences. In our learning scenario the agent-environment interactions consist of agent-data interactions – there is no user simulator as in task-oriented dialogue systems. During each verbal contribution, the DRL agent (1) observes the state of the world via a deep neural network, which models a representation of all sentences raised in the conversation together with a set of candidate responses or agent actions (referred as *clustered actions* in

our approach); (2) it then selects an action so that its word-based representation is sent to the environment; and (3) it receives an updated dialogue history and a numerical reward for having chosen each action, until a termination condition is met. This process—illustrated in Figure 1—is carried out iteratively until the end of a dialogue for as many dialogues as necessary, i.e. until there is no further improvement in the agent's performance.

The contributions of this paper are as follows.

- We propose to train chatbots using value-based deep reinforcement learning using action spaces derived from unsupervised clustering, where each action cluster is a representation of a type of meaning (greeting, question around a topic, statements around a topic, etc.).
- We propose a simple though promising reward function. It is based on human-human dialogues and noisy dialogues for learning to rate good vs. bad dialogues. According to an analysis of dialogue reward prediction, dialogues with lengthy dialogue histories (of at least 10 sentences) report strong correlations between true and predicted rewards on test data.
- Our experiments comparing different sentence embedding sizes (100 vs. 300) did not report statistical differences on test data. This means that similar results can be obtained more efficiently with the smaller embedding than the larger one due to less features. In other words, sentence embeddings of 100 dimensions are as good as 300 dimensions but less computationally demanding.
- Last but not least, we found that training chatbots on multiple data splits is crucial for improved performance over training chatbots using the entire training set.

The remainder of the paper describes our proposed approach in more detail and evaluates it using a publicly available dataset of chitchat conversations. Although our learning agents indeed improve their performance over time with dialogues that they get familiarised with, their performance drops with dialogues that the agents are not familiar with. The former is promising and in favour of our proposed approach, and the latter is not, but it is a general problem faced by data-driven chatbots and an interesting avenue for future research.

## II. RELATED WORK

Reinforcement Learning (RL) methods are typically based on value functions or policy search [2], which also applies to deep RL methods. While value functions have been particularly applied to task-oriented dialogue systems [5]–[10], policy search has been particularly applied to open-ended dialogue systems such as (chitchat) chatbots [11]–[15]. This is not surprising given the fact that task-oriented dialogue systems use finite action sets, while chatbot systems use infinite action sets. So far there is a preference for policy search methods for chatbots, but it is not clear whether they should be preferred because they face problems such as local optima rather than global optima, inefficiency and high variance. It is thus that this paper explores the feasibility of value function-based methods for chatbots, which has not been explored before—at least not from the perspective of deriving the action sets automatically as attempted in this paper.

Other closely related methods to deep RL include seq2seq models for dialogue generation [16]–[21]. These methods tend to be data-hungry because they are typically trained with millions of sentences, which imply high computational demands. While they can be used to address the same problem, in this paper we focus our attention on deep RL-based chatbots and leave their comparison or combination as future work. Nonetheless, these related works agree with the fact that evaluation is a difficult part and that there is a need for better evaluation metrics [22]. This is further supported by [23], where they found that metrics such as Bleu and Meteor amongst others do not correlate with human judgments.

With regard to performance metrics, the reward functions used by deep RL dialogue agents are either specified manually depending on the application, or learnt from dialogue data. For example, [11] conceives a reward function that rewards positively sentences that are easy to respond and coherent while penalising repetitiveness. [12] uses an adversatial approach, where the discriminator is trained to score human vs. non-human sentences so that the generator can use such scores during training. [13] trains a reward function from human ratings. All these related works are neural-based, and there is no clear best reward function to use in future (chitchat) chatbots. This motivated us to propose a new metric that is easy to implement, practical due to requiring only data in raw text, and potentially promising as described below.

## III. PROPOSED APPROACH

To explain the proposed learning approach we first describe how to conceive a finite set of dialogue actions from raw text, then we describe how to assign rewards, and finally describe how to bring everything together during policy learning.

### A. Clustered Actions

Actions in reinforcement learning chatbots correspond to sentences, and their size is infinite assuming all possible combinations of words sequences in a given language. This is especially true in the case of open-ended conversations that make use of large vocabularies, as opposed to task-oriented conversations that make use of smaller (restricted) vocabularies. A **clustered action** is a group of sentences sharing a similar or related meaning via *sentence vectors* derived from word embeddings [24], [25]. While there are multiple ways of selecting features for clustering and also multiple clustering algorithms, the following requirements arise for chatbots: (1) unlabelled data due to human-human dialogues in raw text (this makes it difficult to evaluate the goodness of clustering features and algorithms), and (2) scalability to clustering a large set of data points (sentences in our case, which are mostly unique).

Given a set of data points $\{\mathbf{x}_1, \cdots, \mathbf{x}_n\} \forall \mathbf{x}_i \in \mathbb{R}^m$ and a similarity metric $d(\mathbf{x}_i, \mathbf{x}_{i'})$, the task is to find a set of $k$ clusters with a clustering algorithm. Since in our case each data point

**x** corresponds to a sentence within a dialogue, we represent sentences via their mean word vectors—similarly as in Deep Averaging Networks [26]—denoted as

$$\mathbf{x}_i = \frac{1}{N_i} \sum_{j=1}^{N_i} c_j,$$

where $c_j$ is the vector of coefficients of word $j$ and $N_i$ is the number of words in sentence $i$. For scalability purposes, we use the K-Means++ algorithm [27] with the Euclidean distance

$$d(\mathbf{x}_i^j, \mathbf{x}_{i'}^j) = \sqrt{\sum_{j=1}^{m} (\mathbf{x}_i^j, \mathbf{x}_{i'}^j)^2}$$

with $m$ dimensions, and assume that $k$ is provided rather than automatically induced – though other algorithms can be used with our approach. In this way, a trained clustering model assigns a cluster ID $a \in A$ to features $\mathbf{x}_i$, where the number of actions is equivalent to the number of clusters, i.e. $|A| = k$.

### B. Human-Likeness Rewards

Reward functions in reinforcement learning dialogue agents is often a difficult aspect. We propose to derive the rewards from human-human dialogues by assigning positive values to contextualised responses seen in the data, and negative values to randomly generated responses due to lacking coherence (also referred to as 'non-human-like responses') – see example in Table V. Thus, an episode or dialogue reward can be computed as $R_i = \sum_{j=1}^{N} r_j^i(a)$, where $i$ is the dialogue in focus, $j$ the dialogue turn in focus, and $r_j^i(a)$ is given according to

$$r_j^i(a) = \begin{cases} +1, & \text{if } a \text{ is a human response in dialogue-turn } i, j. \\ -1, & \text{if } a \text{ is human but randomly chosen (incoherent).} \end{cases}$$

### C. Policy Learning

Our Deep Reinforcement Learning (DRL) agents aim to maximise their cumulative reward overtime according to

$$Q^*(s, a; \theta) = \max_{\pi_\theta} \mathbb{E}[r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \cdots | s, a, \pi_\theta],$$

where $r$ is the numerical reward given at time step $t$ for choosing action $a$ in state $s$, $\gamma$ is a discounting factor, and $Q^*(s, a; \theta)$ is the optimal action-value function using weights $\theta$ in a neural network. During training, a DRL agent will choose actions in a probabilistic manner in order to explore new $(s, a)$ pairs for discovering better rewards or to exploit already learnt values—with a reduced level of exploration overtime and an increased level of exploitation over time. During testing, a DRL agent will choose the best actions $a^*$ according to

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

Our DRL agents implement the procedure above using a generalisation of the DQN method [28]—see Algorithm 1. After initialising replay memory $D$, dialogue history $H$, action-value function $Q$ and target action-value function $\hat{Q}$,

we sample a training dialogue from our data of human-human conversations (lines 1-4). A human starts the conversation, which is mapped to its corresponding sentence embedding representation (lines 5-6). Then a set of candidate responses is generated including (1) the true human response and (2) randomly chosen responses (distractors). The candidate responses are clustered as described in Section III-A and the resulting actions are taken into account by the agent for action selection (lines 8-10). Once an action is chosen, it is conveyed to the environment, a reward is observed as described in Section III-B, and the agent's partner response is observed as well in order to update the dialogue history $H$ (lines 11-14). With such an update, the new sentence embedding representation is generated from $H$ in order to update the replay memory $D$ with learning experience $(s, a, r, s')$ (lines 15-16). Then a minibatch of experiences $MB = (s_j, a_j, r_j, s'_j)$ is sampled from $D$ in order to update the weights $\theta$ according to the error derived from the difference between the target value $y_j$ and the predicted value $Q(s, a; \theta)$ (see lines 18 and 20), which is based on the following loss function:

$$L(\theta_j) = \mathbb{E}_{MB} \left( r + \gamma \max_{a'} \hat{Q}(s', a'; \hat{\theta}_j) - Q(s, a; \theta_j) \right)^2.$$

The target action-value function $\hat{Q}$ and state $s$ are updated accordingly (lines 21-22), and this iterative procedure continues until convergence.

## IV. EXPERIMENTS AND RESULTS

### A. Data

We used data from the *Persona-Chat* data set[1], which includes 17,877 dialogues for training (131,431 turns) and 999 dialogues for testing (7,793 turns). They represent averages of 7.35 and 7.8 dialogue turns for training and testing, respectively—see example dialogue in Table V. The vocabulary size in the entire data set contains 19,667 unique words.

### B. Experimental Setting

To analyse the performance of our ChatDQN agents we use subsets of training data vs. the entire training data set. The former are automatically generated by using sentence vectors to represent the features of each dialogue—as described in Section III-A. Similarly, the agents' states are modelled using sentence vectors of the dialogue history with the pretrained coefficients of the Glove model [25]. In all our experiments we use the following neural network architecture[2]:

- mean word vectors, one per sentence, in the input layer (maximum number of vectors=50, with zero-padding) – each word vector of 100 or 300 embedding size,
- two Gated Recurrent Unit (GRU) [30] layers with latent dimensionality of 256, and

| Human Sentences | Distorted Human Sentences |
|---|---|
| hello what are doing today? | hello what are doing today? |
| i'm good, i just got off work and tired, i have two jobs.[r=+1] | do your cats like candy?[r=-1] |
| i just got done watching a horror movie | i just got done watching a horror movie |
| i rather read, i have read about 20 books this year.[r=+1] | do you have any hobbies?[r=-1] |
| wow! i do love a good horror movie. loving this cooler weather | wow! i do love a good horror movie. loving this cooler weather |
| but a good movie is always good.[r=+1] | good job! if you live to 100 like me, you will need all that learning.[r=-1] |
| yes! my son is in junior high and i just started letting him watch them | yes! my son is in junior high and i just started letting him watch them |
| i work in the movies as well.[r=+1] | what a nice gesture. i take my dog to compete in agility classes.[r=-1] |
| neat!! i used to work in the human services field | neat!! i used to work in the human services field |
| yes it is neat, i stunt double, it is so much fun and hard work.[r=+1] | you work very hard. i would like to do a handstand. can you teach it?[r=-1] |
| yes i bet you can get hurt. my wife works and i stay at home | yes i bet you can get hurt. my wife works and i stay at home |
| nice, i only have one parent so now i help out my mom.[r=+1] | yes i do, red is one of my favorite colors[r=-1] |
| i bet she appreciates that very much. | i bet she appreciates that very much. |
| she raised me right, i'm just like her.[r=+1] | haha, it is definitely attention grabbing![r=-1] |
| my dad was always busy working at home depot | my dad was always busy working at home depot |
| now that i am older home depot is my toy r us.[r=+1] | i bet there will be time to figure it out. what are your interests?[r=-1] |

TABLE I

MODIFIED DIALOGUE FROM THE PERSONA-CHAT DATASET [21] WITH OUR PROPOSED REWARDS: $r$=+1 MEANS A HUMAN-LIKE SENTENCE AND $r$=-1 MEANS NON-HUMAN LIKE. THE LATTER SENTENCES, IN RED, ARE SAMPLED RANDOMLY FROM DIFFERENT DIALOGUES IN THE SAME DATASET

---

**Algorithm 1** ChatDQN Learning

1: Initialise Deep Q-Networks with replay memory $D$, dialogue history $H$, action-value function $Q$ with random weights $\theta$, and target action-value functions $\hat{Q}$ with $\hat{\theta} = \theta$
2: Initialise clustering model from training dialogue data
3: **repeat**
4:     Sample a training dialogue (human-human sentences)
5:     Append first sentence to dialogue history $H$
6:     $s$ = sentence embedding representation of $H$
7:     **repeat**
8:         Generate *noisy* candidate response sentences
9:         $A$ = cluster IDs of candidate response sentences
10:         $a = \begin{cases} rand_{a \in A} & \text{if random number} \leq \epsilon \\ \max_{a \in A} Q(s, a; \theta) & \text{otherwise} \end{cases}$
11:         Execute chosen clustered action $a$
12:         Observe human-likeness dialogue reward $r$
13:         Observe environment response (agent's partner)
14:         Append agent and environment responses to $H$
15:         $s'$ = sentence embedding representation of $H$
16:         Append transition $(s, a, r, s')$ to $D$
17:         Sample random minibatch $(s_j, a_j, r_j, s'_j)$ from $D$
18:         $y_j = \begin{cases} r_j & \text{if final step of episode} \\ r_j + \gamma \max_{a' \in A} \hat{Q}(s', a'; \hat{\theta}) & \text{otherwise} \end{cases}$
19:         Set $err = (y_j - Q(s, a; \theta))^2$
20:         Gradient descent step on $err$ with respect to $\theta$
21:         Reset $\hat{Q} = Q$ every $C$ steps
22:         $s \leftarrow s'$
23:     **until** end of dialogue
24:     Reset dialogue history $H$
25: **until** convergence

---

- fully connected layer with number of nodes=the number of clusters, i.e. each cluster corresponding to one action.

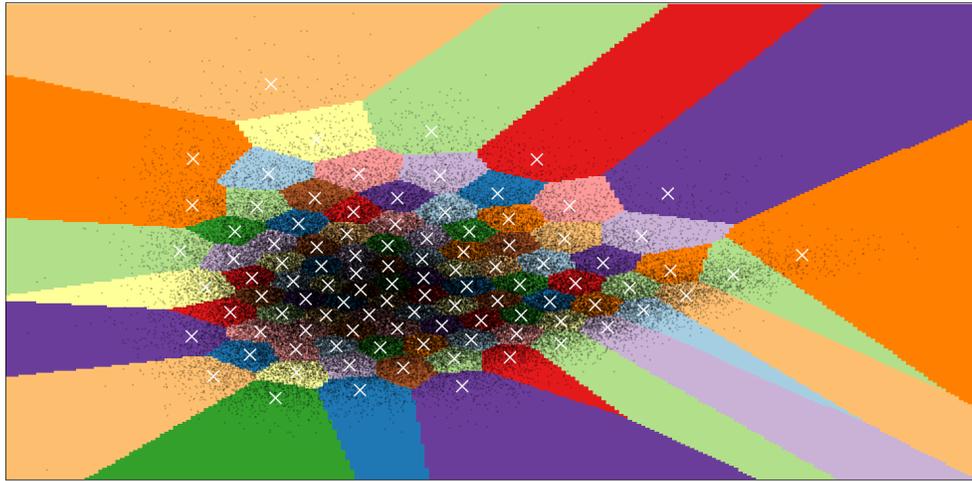While a small number of sentence clusters could result in actions being assigned to potentially the same cluster, a larger number of sentence clusters would mitigate the problem, but the larger the number of clusters the larger the computational expense—i.e. more parameters in the neural network. Figure 2(a) shows an example of our sentence clustering using 100 clusters on our training data. A manual inspection showed that greeting sentences were mostly assigned to the same cluster, and questions expressing preferences (e.g. What is your favourite X?) were also assigned to the same cluster. In this work we thus use a sentence clustering model with $k$=100 derived from our training data and prior to reinforcement learning[3]. In addition, we trained a second clustering model to analyse our experiments using different data splits, where instead of clustering sentences we cluster dialogues. Given that we represent a sentence using a mean word vector, a dialogue can thus be represented by a group of sentence vectors. Figure 2(b) shows an example of our dialogue clustering using 20 clusters on our training data.

Notice that while previous related works in task-oriented DRL-based agents typically use a user simulator, this paper does not use a simulator. Instead, we use the dataset of human-human dialogues directly and substitute one partner conversant in the dialogues by a DRL agent. The goal of the agent is to choose the human-generated sentences (actions) out of a set of candidate responses.
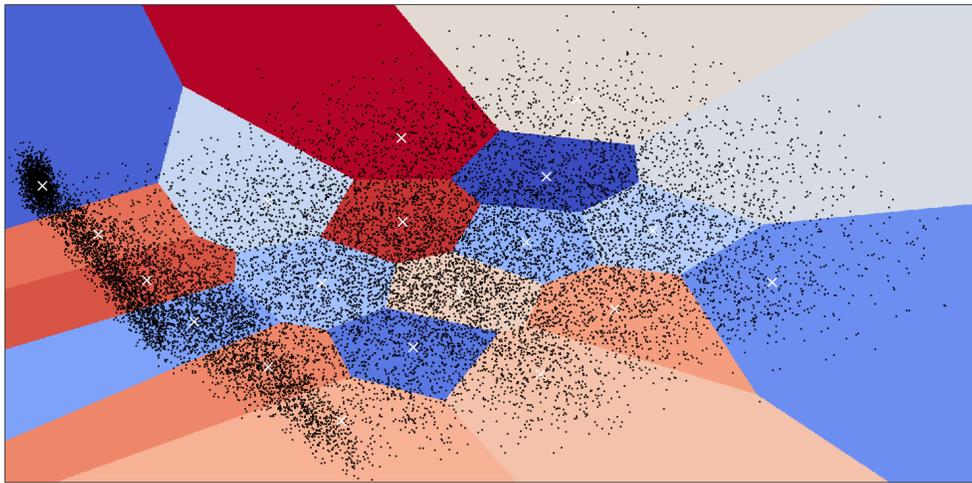
*C. Experimental Results*

The plots in Figure 3 show the training performance of our ChatDQN agents—all using 100 clustered actions. Each plot contains two learning curves, one per agent, where each agent uses a different sentence embedding size (100 or 300 dimensions). In addition, each plot uses an automatically generated data split according to our clustered dialogues. These plots show evidence that all agents indeed improve their behaviour over time even when they use only 100 actions. This can be observed from their average episode rewards, the higher

---

[3]Each experiment in this paper was ran on a GPU Tesla K80 using the following libraries: Keras (https://github.com/keras-team/keras), OpenAI (https://github.com/openai) and Keras-RL (https://github.com/keras-rl/keras-rl).

(a) 100 clusters of training sentences



(b) 20 clusters of training dialogues

Fig. 2. Example clusters of our training data using Principal Component Analysis [31] for visualisations in 2D – black dots represent sentences or dialogues

the better in all learning curves. From a visual inspection, we can observe that the agents using either embedding size (100 or 300) perform rather equivalently but with a small trend for 300 dimensions to dominate its counterpart – more on this below.

The performance of our ChatDQN agents using all training dialogues is shown in Figure 4. It can be noted that in contrast to the previous agents where their improvement in average reward reached values of around 2, the performance in these agents was lower (with average episode reward $< 0$). We attribute this to the larger amount of variation exhibited from about 1K dialogues to 17.8K dialogues.

We analysed the performance of our agents further by using a test set of totally unseen dialogues during training. Table II summarises our results, where we can note that the larger sentence embedding size (300) generally performed better. While a significant difference (according to a two-tailed Wilcoxon Singed Rank Test) at $p = 0.05$ was identified in testing on the training set, *no significant difference* was found

in performance during testing on the test set. These results could be confirmed in other datasets and/or settings in future work. In addition, we can observe that the ChatDQN agents trained using all data (agents with id=20) were not able to achieve as good performance than those agents using smaller data splits. Our results thus reveal that training chatbots on some sort of domains (groups of dialogues automatically discovered in our case), is useful for improved performance.

## V. ANALYSIS OF HUMAN-LIKENESS REWARDS

We employ the algorithm of [32] for extending a dataset of human-human dialogues with distorted dialogues. The latter include varying amounts of distortions, i.e. different degrees of human-likeness. We use such data for training and testing reward prediction models in order to analise the goodness of our proposed reward function. Given extended dataset $\hat{\mathcal{D}} = \{(\hat{d}_1, y_1), \ldots, (\hat{d}_N, y_N)\}$ with (noisy) dialogue histories $\hat{d}_i$, the goal is to predict dialogue scores $y_i$ as accurately as possible. We represent a dialogue history via its sentence

(a) ChatDQN agents using data splits 0 to 4 (from left to right)



(b) ChatDQN agents using data splits 5 to 9 (from left to right)



(c) ChatDQN agents using data splits 10 to 14 (from left to right)



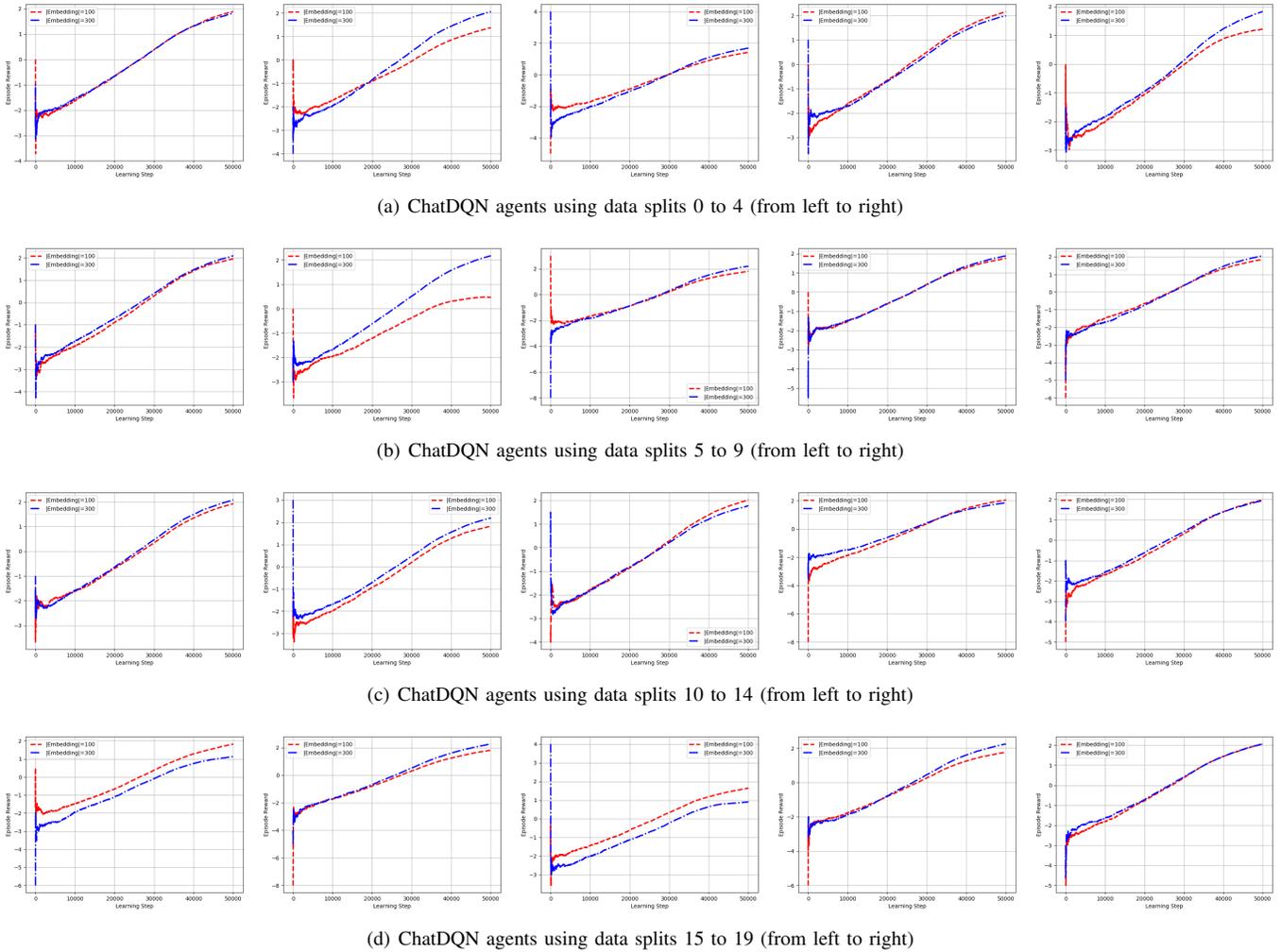(d) ChatDQN agents using data splits 15 to 19 (from left to right)

Fig. 3. Training performance of ChatDQN agents using different data splits of dialogues—see text for details
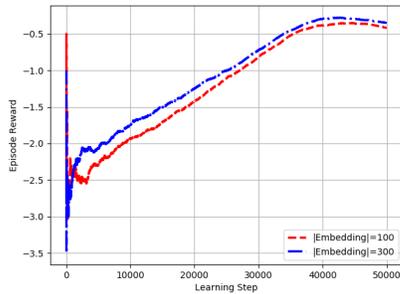


Fig. 4. Training performance of our ChatDQN agents using all training dialogues and two sentence embedding sizes

vectors as in Deep Averaging Networks [26], where sentences are represented with numerical feature vectors denoted as $\mathbf{x} = \{x_1, ..., x_{|\mathbf{x}|}\}$. In this way, a set of word sequences $s_j^i$ in dialogue-sentence pair $i, j$ is mapped to feature vectors

$$\mathbf{x}_j^i = \frac{1}{N_j^i} \sum_{k=1}^{N_j^i} c_{j,k}^i,$$

where $c_{j,k}^i$ is the vector of coefficients of word $k$, part of sentence $j$ in dialogue $i$, and $N_j^i$ is the number of words in the sentence in focus.

Assuming that vector $\mathbf{Y} = \{y_1, ..., y_{|\mathbf{Y}|}\}$ is the set of target labels—generated as described in the dialogue generation algorithm of [32], and using the same test data as the previous section. In this way, dataset $\mathcal{D}^{train} = (\mathbf{X}^{train}, \mathbf{Y}^{train})$ is used for training neural regression models using varying amounts of dialogue history, and dataset $\mathcal{D}^{test} = (\mathbf{X}^{test}, \mathbf{Y}^{test})$ is used for testing the learnt models.

Our experiments use a 2-layer Gated Recurrent Unit (GRU) neural network [30], similar to the one in SectionIV-B but including Batch Normalisation [33] between hidden layers.

We trained neural networks for six different lengths of dialogue history, ranging from 1 sentence to 50 sentences. Each length size involved a separate neural network, trained 10 times in order to report results over multiple runs. Figure 5 reports the average Pearson correlation coefficient—between true dialogue rewards and predicted dialogue rewards—for each length size. It can be observed that short dialogue histories contribute to obtain weak correlations, and that longer

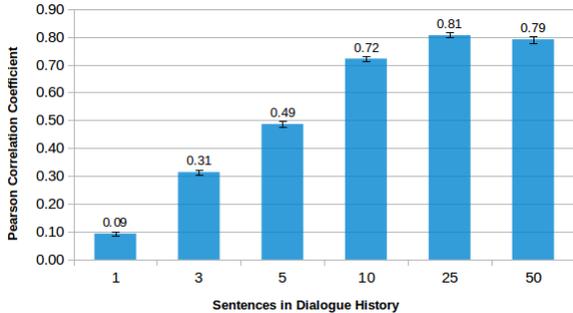| |Embedding|=100 | | | | |Embedding|=300 | | | |
|---|---|---|---|---|---|---|---|
| Data Split (|dialogues|) | Training | Testing on the Training Set | Testing on the Test Set | Data Split (|dialogues|) | Training | Testing on the Training Set | Testing on the Test Set |
| 0 (861) | 1.8778 | 3.7711 | -1.1708 | 0 (1000) | 1.8168 | 3.6785 | -0.8618 |
| 1 (902) | 1.3751 | 3.1663 | -1.7006 | 1 (850) | 2.0622 | 4.4598 | -1.8688 |
| 2 (907) | 1.4194 | 3.1579 | -0.9723 | 2 (1010) | 1.6896 | 3.6724 | -1.4282 |
| 3 (785) | 2.1532 | 4.2508 | -1.3444 | 3 (1029) | 1.9845 | 4.0136 | -0.6109 |
| 4 (1046) | 1.2204 | 2.1581 | -1.5633 | 4 (951) | 1.8255 | 4.0423 | -1.4448 |
| 5 (767) | 1.9456 | 3.9017 | -1.2123 | 5 (832) | 2.0860 | 4.2182 | -0.8277 |
| 6 (1053) | 0.4621 | 0.1370 | -1.8443 | 6 (815) | 2.1735 | 4.2592 | -1.5193 |
| 7 (968) | 1.8090 | 3.8368 | -1.1137 | 7 (891) | 2.1921 | 4.5799 | -1.4233 |
| 8 (858) | 1.7608 | 3.5531 | -1.6678 | 8 (905) | 1.8835 | 3.8337 | -0.6628 |
| 9 (826) | 1.8431 | 3.6254 | -1.0919 | 9 (892) | 2.0521 | 4.1882 | -1.5267 |
| 10 (818) | 1.9188 | 3.8629 | -0.5394 | 10 (835) | 2.0709 | 4.2852 | -0.8831 |
| 11 (944) | 1.8212 | 3.5724 | -1.7020 | 11 (873) | 2.1902 | 4.4848 | -1.3329 |
| 12 (873) | 2.0195 | 4.1895 | -1.3456 | 12 (948) | 1.7761 | 3.7927 | -1.6167 |
| 13 (895) | 2.0515 | 4.1873 | -1.8034 | 13 (932) | 1.8563 | 3.6208 | -1.5149 |
| 14 (863) | 1.9722 | 4.1479 | -1.3244 | 14 (812) | 1.9486 | 4.0347 | -1.5866 |
| 15 (842) | 1.8214 | 3.8942 | -0.8921 | 15 (880) | 1.1338 | 2.4880 | -1.4084 |
| 16 (837) | 1.8162 | 3.8817 | -1.3784 | 16 (787) | 2.2628 | 4.5583 | -1.4290 |
| 17 (958) | 1.6373 | 3.3373 | -0.7726 | 17 (994) | 0.9038 | 1.5106 | -1.5925 |
| 18 (1012) | 1.7631 | 3.6279 | -1.2690 | 18 (853) | 2.2405 | 4.4716 | -1.4231 |
| 19 (862) | 2.0683 | 4.2026 | -1.5901 | 19 (788) | 2.0686 | 4.2219 | -0.9594 |
| 20 (17877) | -0.4138 | -1.2473 | -1.9684 | 20 (17877) | -0.3516 | -0.3490 | -2.0870 |
| Average$^{0-20}$ | 1.6353 | 3.2959† | -1.3461 | Average$^{0-20}$ | **1.8031** | **3.7174†** | **-1.3337** |
| Sum$^{0-20}$ | 34.3419 | 69.2146 | -28.2674 | Sum$^{0-20}$ | **37.8656** | **78.0653** | **-28.0079** |
| Upper Bound | 7.1810 | 7.1810 | 7.5942 | Upper Bound | 7.1810 | 7.1810 | 7.5942 |
| Lower Bound | -7.2834 | -7.2834 | -7.7276 | Lower Bound | -7.2834 | -7.2834 | -7.7276 |
| Random Sel. | -2.4139 | -2.4139 | -2.5526 | Random Sel. | -2.4139 | -2.4139 | -2.5526 |



Fig. 5. Bar plot showing the performance of our dialogue reward predictors using different amounts of dialogue history (from 1 sentence to 50 sentences). Each bar reports an average Pearson correlation score over 10 runs, where the coefficients report the correlation between true dialogue rewards and predicted dialogue rewards in our test data

without any manual annotations. The task of the agents is to learn to choose human-like actions (sentences) out of candidate responses including human generated and randomly generated sentences. In our proposed rewards we assume that the latter are generally incoherent throughout the dialogue history. Experimental results using chitchat data report that DRL agents learn reasonable policies using training dialogues, but their generalisation ability in a test set of unseen dialogues remains a key challenge for future research in this field. In addition, we found the following: (a) that sentence embedding sizes of 100 and 300 perform equivalently on test data; (b) that training agents using larger amounts of training can deteriorate performance than training with smaller amounts; and (c) that our proposed dialogue rewards can be predicted with strong correlation (between true and predicted rewards) by using neural-based regressors with lengthy dialogue histories of $\geq$ 10 sentences (25 sentences was the best in our experiments).

Future work can explore the following avenues. First, confirm these findings with other datasets and settings in order to draw even stronger conclusions. Second, investigate further the proposed approach for improved generalisation in test data. For example, other methods of feature extraction, clustering algorithms, distance metrics, policy learning algorithms, architectures, and a comparison of reward functions can be explored. Last but not least, combine the proposed learning approach with more knowledge intensive resources [34], [35] such as semantic parsers, coreference resolution, among others.

dialogue histories ($\geq$ 10 sentences) contribute to obtain strong correlations. It can also be observed that the longest history may not be the best choice of length size, the network using 25 sentences achieved the best results. From these results we can conclude that our proposed human-likeness rewards—with lengthy dialogue histories—can be used for training future neural-based chatbots.

## VI. CONCLUSION AND FUTURE WORK

This paper presents a novel approach for training Deep Reinforcement Learning (DRL) chatbots, which uses clustered actions and rewards derived from human-human dialogues

REFERENCES

[1] Barbara J. Grosz and Candace L. Sidner, "Attention, intentions, and the structure of discourse," *Computational Linguistics*, vol. 12, no. 3, pp. 175–204, 1986.

[2] Richard S. Sutton and Andrew G. Barto, *Reinforcement learning - an introduction*, Adaptive computation and machine learning. MIT Press, 2nd edition edition, 2018.

[3] Trevor Hastie, Robert Tibshirani, and Jerome H. Friedman, *The elements of statistical learning: data mining, inference, and prediction, 2nd Edition*, Springer series in statistics. Springer, 2009.

[4] Yann LeCun, Yoshua Bengio, and Geoffrey E. Hinton, "Deep learning," *Nature*, vol. 521, no. 7553, pp. 436–444, 2015.

[5] Iñigo Casanueva, Pawel Budzianowski, Pei-Hao Su, Stefan Ultes, Lina Maria Rojas-Barahona, Bo-Hsiang Tseng, and Milica Gasic, "Feudal reinforcement learning for dialogue management in large domains," in *NAACL-HLT*, 2018.

[6] Heriberto Cuayáhuitl, "SimpleDS: A simple deep reinforcement learning dialogue system," *CoRR*, vol. abs/1601.04574, 2016.

[7] Heriberto Cuayáhuitl, Seunghak Yu, Ashley Williamson, and Jacob Carse, "Scaling up deep reinforcement learning for multi-domain dialogue systems," in *IJCNN*, 2017.

[8] Heriberto Cuayáhuitl and Seunghak Yu, "Deep reinforcement learning of dialogue policies with less weight updates," in *INTERSPEECH*, 2017.

[9] Jason D. Williams, Kavosh Asadi, and Geoffrey Zweig, "Hybrid code networks: practical and efficient end-to-end dialog control with supervised and reinforcement learning," in *ACL*, 2017.

[10] Baolin Peng, Xiujun Li, Lihong Li, Jianfeng Gao, Asli Çelikyilmaz, Sungjin Lee, and Kam-Fai Wong, "Composite task-completion dialogue policy learning via hierarchical deep reinforcement learning," in *EMNLP*, 2017.

[11] Jiwei Li, Will Monroe, Alan Ritter, Dan Jurafsky, Michel Galley, and Jianfeng Gao, "Deep reinforcement learning for dialogue generation," in *EMNLP*, 2016.

[12] Jiwei Li, Will Monroe, Tianlin Shi, Sébastien Jean, Alan Ritter, and Dan Jurafsky, "Adversarial learning for neural dialogue generation," in *EMNLP*, 2017.

[13] Iulian Vlad Serban, Chinnadhurai Sankar, Mathieu Germain, Saizheng Zhang, Zhouhan Lin, Sandeep Subramanian, Taesup Kim, Michael Pieper, Sarath Chandar, Nan Rosemary Ke, Sai Rajeswar, Alexandre de Brébisson, Jose M. R. Sotelo, Dendi Suhubdy, Vincent Michalski, Alexandre Nguyen, Joelle Pineau, and Yoshua Bengio, "A deep reinforcement learning chatbot (short version)," *CoRR*, vol. abs/1801.06700, 2018.

[14] Chinnadhurai Sankar and Sujith Ravi, "Modeling non-goal oriented dialog with discrete attributes," in *NeurIPS Workshop on Conversational AI: "Today's Practice and Tomorrow's Potential"*, 2018.

[15] Chih-Wei Lee, Yau-Shian Wang, Tsung-Yuan Hsu, Kuan-Yu Chen, Hung-yi Lee, and Lin-Shan Lee, "Scalable sentiment for sequence-to-sequence chatbot response with performance analysis," *CoRR*, vol. abs/1804.02504, 2018.

[16] Oriol Vinyals and Quoc V. Le, "A neural conversational model," *CoRR*, vol. abs/1506.05869, 2015.

[17] Alessandro Sordoni, Michel Galley, Michael Auli, Chris Brockett, Yangfeng Ji, Margaret Mitchell, Jian-Yun Nie, Jianfeng Gao, and Bill Dolan, "A neural network approach to context-sensitive generation of conversational responses," in *HLT-NAACL*, 2015.

[18] Iulian Vlad Serban, Tim Klinger, Gerald Tesauro, Kartik Talamadupula, Bowen Zhou, Yoshua Bengio, and Aaron C. Courville, "Multiresolution recurrent neural networks: An application to dialogue response generation," in *AAAI*, 2017.

[19] Jiwei Li, Michel Galley, Chris Brockett, Georgios P. Spithourakis, Jianfeng Gao, and William B. Dolan, "A persona-based neural conversation model," in *ACL*, 2016.

[20] Wenjie Wang, Minlie Huang, Xin-Shun Xu, Fumin Shen, and Liqiang Nie, "Chat more: Deepening and widening the chatting topic via a deep model," in *SIGIR*. 2018, ACM.

[21] Saizheng Zhang, Emily Dinan, Jack Urbanek, Arthur Szlam, Douwe Kiela, and Jason Weston, "Personalizing dialogue agents: I have a dog, do you have pets too?," *CoRR*, vol. abs/1801.07243, 2018.

[22] Rui Yan, ""chitty-chitty-chat bot": Deep learning for conversational AI," in *IJCAI*, 2018.

[23] Chia-Wei Liu, Ryan Lowe, Iulian Serban, Michael Noseworthy, Laurent Charlin, and Joelle Pineau, "How NOT to evaluate your dialogue system: An empirical study of unsupervised evaluation metrics for dialogue response generation," in *EMNLP*, 2016.

[24] Tomas Mikolov, Ilya Sutskever, Kai Chen, Gregory S. Corrado, and Jeffrey Dean, "Distributed representations of words and phrases and their compositionality," in *NIPS*, 2013.

[25] Jeffrey Pennington, Richard Socher, and Christopher D. Manning, "Glove: Global vectors for word representation," in *EMNLP*, 2014.

[26] Mohit Iyyer, Varun Manjunatha, Jordan L. Boyd-Graber, and Hal Daumé III, "Deep unordered composition rivals syntactic methods for text classification," in *ACL (1)*, 2015.

[27] David Arthur and Sergei Vassilvitskii, "K-means++: The advantages of careful seeding," in *SODA*. 2007, SIAM.

[28] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A. Rusu, Joel Veness, Marc G. Bellemare, Alex Graves, Martin A. Riedmiller, Andreas Fidjeland, Georg Ostrovski, Stig Petersen, Charles Beattie, Amir Sadik, Ioannis Antonoglou, Helen King, Dharshan Kumaran, Daan Wierstra, Shane Legg, and Demis Hassabis, "Human-level control through deep reinforcement learning," *Nature*, vol. 518, no. 7540, 2015.

[29] Alexander H. Miller, Will Feng, Dhruv Batra, Antoine Bordes, Adam Fisch, Jiasen Lu, Devi Parikh, and Jason Weston, "Parlai: A dialog research software platform," in *EMNLP (System Demonstrations)*, 2017.

[30] Kyunghyun Cho, Bart van Merrienboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio, "Learning phrase representations using RNN encoder–decoder for statistical machine translation," in *EMNLP*. 2014, Association for Computational Linguistics.

[31] M. E. Tipping and Christopher Bishop, "Probabilistic principal component analysis," *Journal of the Royal Statistical Society, Series B*, vol. 21/3, pp. 611622, January 1999.

[32] Heriberto Cuayáhuitl, Seonghan Ryu, Donghyeon Lee, and Jihie Kim, "A study on dialogue reward prediction for open-ended conversational agents," in *NeurIPS Workshop on Conversational AI: "Today's Practice and Tomorrow's Potential"*, 2018.

[33] Sergey Ioffe and Christian Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," in *International Conference on Machine Learning (ICML)*, 2015.

[34] Christopher D. Manning, Mihai Surdeanu, John Bauer, Jenny Finkel, Steven J. Bethard, and David McClosky, "The Stanford CoreNLP natural language processing toolkit," in *Association for Computational Linguistics (ACL) System Demonstrations*, 2014.

[35] Nina Dethlefs, "Domain transfer for deep natural language generation from abstract meaning representations," *IEEE Comp. Int. Mag.*, vol. 12, no. 3, pp. 18–28, 2017.