

---

# Improving Robustness of Neural Dialog Systems in a Data-Efficient Way with Turn Dropout

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Neural network-based dialog models often lack robustness to anomalous, out-of-  
2 domain (OOD) user input which leads to unexpected dialog behavior and thus  
3 considerably limits such models' usage in mission-critical production environments.  
4 The problem is especially relevant in the setting of dialog system bootstrapping  
5 with limited training data and no access to OOD examples. In this paper, we explore  
6 the problem of robustness of such systems to anomalous input and the associated  
7 to it trade-off in accuracies on seen and unseen data. We present a new dataset for  
8 studying the robustness of dialog systems to OOD input, which is bAbI Dialog  
9 Task 6 augmented with OOD content in a controlled way. We then present turn  
10 dropout, a simple yet efficient negative sampling-based technique for improving  
11 robustness of neural dialog models. We demonstrate its effectiveness applied to  
12 Hybrid Code Networks (HCNs) on our data. Specifically, an HCN trained with  
13 turn dropout achieves more than **75%** per-utterance accuracy on the augmented  
14 dataset's OOD turns and **74%** F1-score as an OOD detector. Furthermore, we  
15 introduce a Variational HCN enhanced with turn dropout which achieves more  
16 than **56.5%** accuracy on the original bAbI Task 6 dataset, thus outperforming the  
17 initially reported HCN's result.

## 18 1 Introduction

19 Data-driven approaches for building dialog systems have recently passed the stage of open-ended  
20 academic research and are adopted in platforms like *Google Dialogflow*, *Apple SiriKit*, *Amazon*  
21 *Alexa Skills Kit*, and *Microsoft Cognitive Services*. However, most of those platforms' data-driven  
22 functionality is limited to Natural Language Understanding: user intent detection, named entity  
23 recognition, and slot filling. A more unified approach to dialog system bootstrapping — end-to-end  
24 dialog learning — is still only emerging as a commercial service, e.g. *Microsoft Conversation Learner*.  
25 Although still in its early age, end-to-end dialog learning from examples offers great potential: it  
26 doesn't require advanced programming skills and thus it makes it possible for a wider range of users  
27 to create dialog systems for their purposes. In turn, in the enterprise environment, end-to-end dialog  
28 learning bridges the gap between user experience designers and the actual working systems thus  
29 making product cycles and overall workflow faster.

30 From the technical point of view, the key issue in end-to-end training is the lack of robustness of  
31 the resulting systems. In the real-world setting of rapid dialog system prototyping, it is common to  
32 have only in-domain (IND) data for a closed target domain. This leads to a significant overfitting  
33 of machine learning methods and unpredictable behavior in the cases outside of what was seen  
34 during training. For a closed-domain dialog system, it's extremely important to maintain predictable  
35 behavior on anomalous, OOD user input.

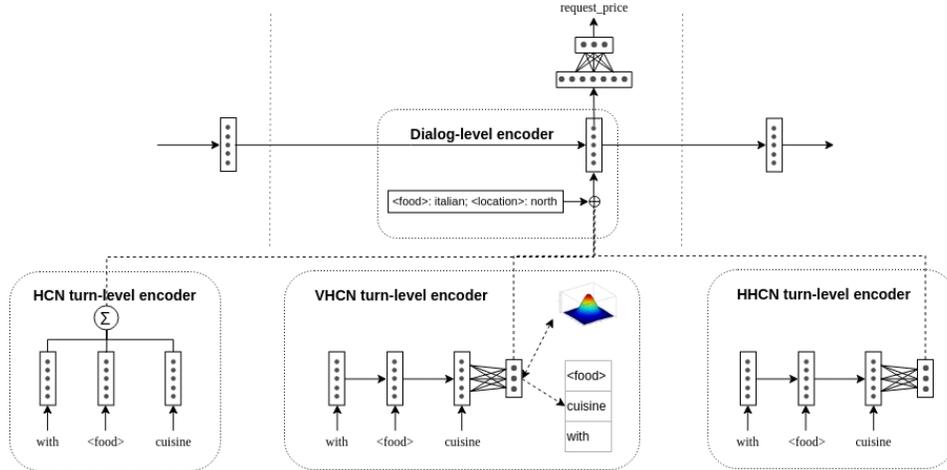


Figure 1: Hybrid Code Network model family

36 In this paper, we focus on studying the effect of OOD input on end-to-end goal-oriented dialog  
 37 models’ performance and propose a simple and efficient solution to improving robustness only using  
 38 IND data. Our contribution is thus two-fold:

- 39 • We present a dataset for studying the effect of OOD input on dialog models.
- 40 • We present turn dropout, an efficient negative sampling-technique for training dialog models  
 41 that are capable of OOD handling using only the IND data.

42 We show that HCN-based models enhanced with turn dropout show superior performance on OOD  
 43 input, as well as surpass original HCN’s result on IND-only data.

## 44 2 Related work

45 Detection of anomalous input is a key research problem in machine learning. In the area of dialog  
 46 systems, there is a series of approaches to detecting and processing of OOD input. If treated as  
 47 a classification problem, this problem require both IND and OOD data [10, 13]. Although for  
 48 the real-world scenario of end-to-end dialog system learning the task of collecting data covering  
 49 potentially unbounded variety of OOD input is impractical. In contrast, there are also approaches like  
 50 an in-domain verification method [7] and an autoencoder-based OOD detection [11] which do not  
 51 require OOD data. However, they still have restrictions such that there must be multiple sub-domains  
 52 to learn utterance representation and one must set a decision threshold for OOD detection. For a  
 53 dialog system that is supposed to work in a single closed domain, these methods are not a viable  
 54 solution.

55 In contrast to those approaches, we present a simple and efficient technique for training dialog systems  
 56 robust to OOD input in an end-to-end way, which allows the model to leverage the dialog context  
 57 information to avoid the necessity of using IND data.

## 58 3 Dataset for studying robustness of dialog systems

59 In order to study the effect of OOD input on end-to-end dialog system’s performance, we used a  
 60 dataset of real human-computer goal-oriented dialogs and augmented it with real user utterances  
 61 from other domains in a controlled way using the open-source toolkit bAbI tools<sup>1</sup> [12].

62 As our main dataset, we use bAbI Dialog Task 6 [2], real human-computer conversations in the  
 63 restaurant serach domain initially collected for Dialog State Tracking Challenge 2 [5].

64 Our OOD augmentations are as follows:

<sup>1</sup>[https://bit.ly/babi\\_tools](https://bit.ly/babi_tools)

- 65 • *turn-level OOD*: user requests from a foreign domain — the desired system behavior for  
66 such input is the fallback action,
- 67 • *segment-level OOD*: interjections in the user in-domain requests — treated as valid user  
68 input and is supposed to be handled by the system in a regular way.

69 These two augmentation types reflect a specific dialog pattern of interest (see Table 2): first, the user  
70 utters a request from another domain at an arbitrary point in the dialog (each turn is augmented with  
71 the probability  $p_{ood\_start}$ ), and the system answers accordingly. This may go on for several turns  
72 in a row —each following turn is augmented with the probability  $p_{ood\_cont}$ . Eventually, the OOD  
73 sequence ends up and the dialog continues as usual, with a segment-level OOD of the user affirming  
74 their mistake. For this study, we set  $p_{ood\_start}$  to 0.2 and  $p_{ood\_cont}$  to 0.4<sup>2</sup>.

75 While we introduce the OOD augmentations in a controlled programmatic way, the actual OOD  
76 content is natural. The turn-level OOD utterances are taken from dialog datasets in several foreign  
77 domains:

- 78 • Frames dataset [1] — travel booking (1198 utterances),
- 79 • Stanford Key-Value Retrieval Network Dataset [4] — calendar scheduling, weather informa-  
80 tion retrieval, city navigation (3030 utterances),
- 81 • Dialog State Tracking Challenge 1 [15] — bus information (968 utterances).

82 In order to avoid incomplete/elliptical phrases, we only took the first user’s utterances from the  
83 dialogs.

84 For segment-level OOD, we mined utterances with the explicit affirmation of a mistake from Twitter  
85 and Reddit conversations datasets — 701 and 500 utterances respectively. Our datasets, as well as the  
86 tools for OOD-augmentation of arbitrary datasets of interest are openly available<sup>3</sup>.

## 87 4 A data-efficient technique for training robust dialogue systems

### 88 4.1 Models

89 In this paper, we experiment with Hybrid Code Network family of models [14]. HCN is reported  
90 to be state-of-the-art for the original, IND-only bAbI Dialog Task 6 data. Thus, in this paper we  
91 experiment with it and explore its robustness to OOD input.

92 HCN is a hierarchical dialog control model with a turn-level and a dialog-level components (we will  
93 call them both encoders). The turn-level encoder produces a latent representation of a single dialog  
94 turn, and the dialog-level one augments it with additional dialog-level features and produces a latent  
95 representation of the entire dialog with an RNN. The resulting representation is then fed into the  
96 predictor MLP which outputs the final dialog actions (restricted by expert-provided binary action  
97 masks). Our models are described below — they share the same dialog-level encoder and predictor.  
98 The differences are on the turn level and in the overall optimization objective (see Figure 1 for an  
99 illustration).

**HCN** — the original model introduced by [14]. Its encoding of the turn  $t$  of  $N$  tokens is as follows:

$$HCN_t = \frac{1}{N} \sum_i word2vec(t_i)$$

100 Where *word2vec* is the pre-trained Google News word2vec embeddings (frozen at the training time).  
101 HCN’s optimization objective is categorical cross-entropy with respect to negative log-likelihood  
102 (NLL) of resulting actions.

**Hierarchical HCN (HHCN)** uses an RNN (in our case an LSTM cell [6]) for encoding each  
utterance:

$$HHCN_t = LSTM(t)$$

<sup>2</sup>We experimented with other values of  $p_{ood\_start}$  and  $p_{ood\_cont}$  but didn’t see significant differences in the  
results. Further experiments for different domains are encouraged using the tools provided

<sup>3</sup>See <link anonymized>

103 The optimization objective is the same as of HCN. Variants of this model were described by [8] and  
104 [9].

**Variational HCN (VHCN)** which, to the best of our knowledge, is presented here for the first time — uses a Variational Autoencoder as the turn-level encoder, so that the resulting turn encoding is VAE’s latent variable (normally referred to as  $z$ ):

$$VHCN_t = \mu(LSTM(t)) + \sigma(LSTM(t)) * N(0, 1)$$

105 Where  $\mu$  and  $\sigma$  are MLPs for predicting  $z$ ’s posterior distribution parameters, and  $N(0, 1)$  is a sample  
106 from its prior distribution, a standard Gaussian [3].

This model differs from the previous two in that it learns dialog control and autoencoding jointly. In order to keep the secondary task less complex than the main one, we represent VAE’s reconstruction targets as bags of words (BoW). Thus, VHCN loss function is as follows:

$$L_{VHCN} = L_{NLL(a_i)} + L_{BoW(t_i)} + L_{KL(q_z||p_z)}$$

107 In the above formula, the first term is the main task’s NLL loss for the dialog action  $a_i$ , the second  
108 one is the VAE’s BoW reconstruction loss for the input turn  $t_i$ , and the last term is  $KL$ -divergence  
109 between the prior and posterior distribution of the VAE’s latent variable  $z$  — following [3], we  
110 compute it in a closed form.

111 Another benefit of the BoW loss is, as reported in [16], it helps keep the variational properties of the  
112 model (i.e. non-zero KL-term) without the necessity of using the KL-term annealing trick [3] which  
113 is itself challenging to control in practice. Unlike the authors of the original BoW loss approach, we  
114 don’t stack softmax cross-entropy losses for each token and instead use a single sigmoid cross-entropy  
115 loss for the entire BoW vector.

116 All the models above use the same dialog-level LSTM encoder with additional features concatenated  
117 to the turn representations: BoW turn features, dialog context features, and previous system action<sup>4</sup>.

## 118 4.2 Turn dropout

119 In order to train a system robust to OOD in the absence of real OOD examples, we employ a negative  
120 sampling-based approach and generate them synthetically from available IND data with a technique  
121 we call *turn dropout*. Namely, we replace random dialog turns with synthetic ones, and assign them  
122 the fallback action.

123 More formally, our dialog features are as follows:  $\langle f\_turn, f\_ctx, f\_mask, a \rangle$ , i.e. turn features  
124 (token sequences), dialog context features, action masks, and target actions respectively.

125 Under turn dropout, for a randomly selected dialog  $i$  and its turn  $j$ , we replace  $f\_turn_{ij}$  with a  
126 sequence of random vocabulary words (drawn from a uniform distribution over the vocabulary) and  
127 UNK tokens, and corresponding  $a_{ij}$  with the fallback action, and leave all other features intact. In  
128 this way, we’re simulating anomalous turns for the system given usual contexts (as stored in the  
129 dialog RNN’s state), and we put minimum assumptions on the synthesized turns’ structure (we only  
130 limit their lengths to be within the bounds of the real utterances).

## 131 5 Experimental setup and evaluation

132 We train our models only using the original bAbI Dialog Task 6 dataset, and evaluate them on our  
133 OOD-augmented versions of it: we use the per-utterance accuracy as our main evaluation metric;  
134 the models are trained with the same hyperparameters (where applicable) listed in Table 3. The  
135 models use the common unified vocabulary including all words from our datasets (including OOD  
136 content): the intuition behind this is as follows: production dialog models often use word embedding  
137 matrices with vocabularies significantly exceeding that of the training data in order to take advantage  
138 of additional generalization power via relations like synonymy, hyponymy, or hypernymy normally  
139 efficiently handled by distributed word representations. Therefore, mapping every unseen word to an  
140 ‘UNK’ doesn’t quite reflect that setting.

---

<sup>4</sup>Without the loss of the architecture generality, we have action mask vectors as additional features for the dialog-level LSTM [14], but they don’t convey any information and are always set to 1’s

Model	bAbI Dialog Task 6	bAbI Dialog Task 6 + OOD			
	Overall acc.	Overall acc.	Seg. OOD acc.	OOD acc.	OOD F1
HCN	0.557	0.438	<b>0.455</b>	0.0	0.0
HHCN	0.531	0.418	0.424	0.0	0.0
VHCN	0.533	0.413	0.413	0.0	0.0
TD-HCN	0.563	<b>0.575</b>	0.257	<b>0.754</b>	<b>0.743</b>
TD-HHCN	0.505	0.455	0.435	0.274	0.418
TD-VHCN	<b>0.565</b>	0.545	0.407	0.530	0.667

Table 1: Evaluation results

141 We tuned our models’ hyperparameters using 2-stage grid search, tracking the development set  
142 accuracy. At the first stage, we adjusted the embedding dimensionality of our models (and the latent  
143 variable size in case of VHCN). Then, given the values found, at the second stage we adjusted turn  
144 dropout ratio at the interval  $[0.05 - 0.7]$ . Exact hyperparameter values are detailed in Table 3.

145 The results are shown in Table 1 — please note, apart from the accuracies we report OOD F1-measure,  
146 a metric showing the model’s performance as a conventional OOD detector, with positive class being  
147 the fallback action, and negative — all the IND classes actions.

148 Finally, given the stochastic nature of VHCN, we reported its mean accuracy scores over 3 runs (we  
149 used the same criterion for selecting the best model during the training procedure).

## 150 6 Discussion and future work

151 In this paper, we explored the problem of robustness of neural dialog systems to OOD input. Specifi-  
152 cally, we presented a dataset for studying this problem along with a general procedure for augmenting  
153 arbitrary datasets of interest for such purpose. Secondly, we introduced turn dropout, a simple yet  
154 efficient technique for improving OOD robustness of dialog control models and evaluated its effect  
155 on several Hybrid Code Network-family models.

156 As our experiments showed, while learning to handle both IND and OOD input with access to  
157 IND-only data at the training time, there appears the following trade-off: a model performing better  
158 on the ‘clean’ test turns is prone to lower accuracy on OOD — it can be said that it slightly overfits  
159 to its devset. On the other hand, a model regularized with turn dropout during training naturally  
160 performs better on unseen OOD turns, but with not as high accuracy on its ‘clean’, IND test data.  
161 Another side of the trade-off is the accuracy of OOD detection vs robust handling of IND input  
162 with segment-level noise. As our results showed, models specifically trained for OOD detection all  
163 demonstrate lower accuracy on the noisy IND.

164 Among the models we evaluated, it’s worth noting that the original HCN demonstrated the best  
165 performance as an OOD detector (more than **74%** F1-score) and thus overall IND + OOD accuracy  
166 on the augmented dataset — more than **57%**. While some parts of its architecture (e.g. mean  
167 vector-based turn encoding or bag-of-words feature vector at the utterance level) may not seem to be  
168 the most robust solution, the model demonstrate superior overall performance. Averaging at the turn  
169 level instead of recurrent encoding (the case of HHCN and VHCN) makes the model less dependent  
170 on actual word sequences seen during training but on the keywords themselves.

171 In turn, VHCN demonstrated superior performance on IND data when trained with turn dropout,  
172 more than **56%** — it benefited in terms of both overall accuracy and the absence of false-positive  
173 OODs thus outperforming the original HCN as reported by [14]. An additional challenge was to train  
174 it while keeping its variational properties (i.e. reasonably high KL term) — the BoW reconstruction  
175 loss which we used in order to simplify the secondary task, helped with this as well [16]. On the  
176 other hand, while achieving superior performance on clean data, VHCN’s properties didn’t result in  
177 OOD handling improvements.

178 The question which is still unanswered is how these techniques apply to the setting of few-shot  
179 training. In the practical setup of training dialog systems from minimal data, having access to even  
180 medium-sized datasets like bAbI Dialog Task 6 isn’t realistic, and all the initial requirements for  
181 the models have to be met only using the minimal training data available. It’s the next step in our

182 research to explore how our techniques apply to this setup and what needs to be done in order to  
183 achieve OOD robustness with maximum few-shot data efficiency.

## 184 References

- 185 [1] Layla El Asri, Hannes Schulz, Shikhar Sharma, Jeremie Zumer, Justin Harris, Emery Fine,  
186 Rahul Mehrotra, and Kaheer Suleman. Frames: a corpus for adding memory to goal-oriented  
187 dialogue systems. In *Proceedings of the 18th Annual SIGdial Meeting on Discourse and*  
188 *Dialogue, Saarbrücken, Germany, August 15-17, 2017*, pages 207–219, 2017.
- 189 [2] Antoine Bordes, Y-Lan Boureau, and Jason Weston. Learning end-to-end goal-oriented dialog.  
190 *ICLR*, 2017.
- 191 [3] Samuel R. Bowman, Luke Vilnis, Oriol Vinyals, Andrew M. Dai, Rafal Józefowicz, and Samy  
192 Bengio. Generating sentences from a continuous space. In *Proceedings of the 20th SIGNLL*  
193 *Conference on Computational Natural Language Learning, CoNLL 2016, Berlin, Germany,*  
194 *August 11-12, 2016*, pages 10–21, 2016.
- 195 [4] Mihail Eric, Lakshmi Krishnan, Francois Charette, and Christopher D. Manning. Key-value  
196 retrieval networks for task-oriented dialogue. In *Proceedings of the 18th Annual SIGdial*  
197 *Meeting on Discourse and Dialogue, Saarbrücken, Germany, August 15-17, 2017*, pages 37–49,  
198 2017.
- 199 [5] Matthew Henderson, Blaise Thomson, and Jason D. Williams. The second dialog state tracking  
200 challenge. In *Proceedings of the SIGDIAL 2014 Conference, The 15th Annual Meeting of the*  
201 *Special Interest Group on Discourse and Dialogue, 18-20 June 2014, Philadelphia, PA, USA,*  
202 *pages 263–272, 2014.*
- 203 [6] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Comput.*, 9(8):1735–  
204 1780, November 1997.
- 205 [7] Ian Lane, Tatsuya Kawahara, Tomoko Matsui, and Satoshi Nakamura. Out-of-domain utterance  
206 detection using classification confidences of multiple topics. *IEEE Transactions on Audio,*  
207 *Speech, and Language Processing*, 15(1):150–161, 2007.
- 208 [8] Sungjin Lee. Toward continual learning for conversational agents. *CoRR*, abs/1712.09943,  
209 2017.
- 210 [9] Weiri Liang and Meng Yang. Hierarchical hybrid code networks for task-oriented dialogue.  
211 In De-Shuang Huang, Kang-Hyun Jo, and Xiao-Long Zhang, editors, *Intelligent Computing*  
212 *Theories and Application*, pages 194–204, Cham, 2018. Springer International Publishing.
- 213 [10] Mikio Nakano, Shun Sato, Kazunori Komatani, Kyoko Matsuyama, Kotaro Funakoshi, and  
214 Hiroshi G Okuno. A two-stage domain selection framework for extensible multi-domain spoken  
215 dialogue systems. In *Proceedings of the SIGDIAL 2011 Conference*, pages 18–29. Association  
216 for Computational Linguistics, 2011.
- 217 [11] Seonghan Ryu, Seokhwan Kim, Junhwi Choi, Hwanjo Yu, and Gary Geunbae Lee. Neural  
218 sentence embedding using only in-domain sentences for out-of-domain sentence detection in  
219 dialog systems. *Pattern Recognition Letters*, 88:26–32, 2017.
- 220 [12] Igor Shalymov, Arash Eshghi, and Oliver Lemon. Challenging Neural Dialogue Models with  
221 Natural Data: Memory Networks Fail on Incremental Phenomena. In *Proceedings of the 21st*  
222 *Workshop on the Semantics and Pragmatics of Dialogue (SemDial 2017 - SaarDial)*, 2017.
- 223 [13] Gokhan Tur, Anoop Deoras, and Dilek Hakkani-Tür. Detecting out-of-domain utterances  
224 addressed to a virtual personal assistant. In *Fifteenth Annual Conference of the International*  
225 *Speech Communication Association*, 2014.
- 226 [14] Jason D. Williams, Kavosh Asadi, and Geoffrey Zweig. Hybrid code networks: practical and  
227 efficient end-to-end dialog control with supervised and reinforcement learning. In *Proceedings*  
228 *of the 55th Annual Meeting of the Association for Computational Linguistics, ACL 2017,*  
229 *Vancouver, Canada, July 30 - August 4, Volume 1: Long Papers*, pages 665–677, 2017.

- 230 [15] Jason D. Williams, Antoine Raux, Deepak Ramachandran, and Alan W. Black. The dialog state  
 231 tracking challenge. In *Proceedings of the SIGDIAL 2013 Conference, The 14th Annual Meeting*  
 232 *of the Special Interest Group on Discourse and Dialogue, 22-24 August 2013, SUPELEC, Metz,*  
 233 *France*, pages 404–413, 2013.
- 234 [16] Tiancheng Zhao, Ran Zhao, and Maxine Eskénazi. Learning discourse-level diversity for neural  
 235 dialog models using conditional variational autoencoders. In *Proceedings of the 55th Annual*  
 236 *Meeting of the Association for Computational Linguistics, ACL 2017, Vancouver, Canada, July*  
 237 *30 - August 4, Volume 1: Long Papers*, pages 654–664, 2017.

238 **Appendix A**

1 sys: Hello, welcome to the Cambridge restaurant system. You can ask for restaurants by area,  
 price range or food type. How may I help you?  
 2 usr: **I'd like to get away from Monday, August 15, 2016 to Wednesday, August 31, 2016.**  
**My budget is \$3200 and I would leave from Detroit.**  
 3 sys: Sorry I didn't catch that. Could you please repeat?  
 4 usr: **Will there be frost in Durham next week?**  
 3 sys: Sorry I didn't catch that. Could you please repeat?  
 4 usr: **So sorry man** I want a moderately priced restaurant in the west part of town  
 5 sys: <api call>

Table 2: Augmented dialog example (OOD content in bold)

Hyperparameter	HCN	HHCN	VHCN
Embedding size	64	128	128
Latent variable size	—	—	8
Learning rate	0.001		
Optimizer	Adam		
Early stopping threshold (epochs)	20		
Turn dropout ratio	0.4	0.6	0.3
Word dropout ratio	0.2		

Table 3: Model hyperparameters