

Overview of the seventh Dialog System Technology Challenge: DSTC7

D'Haro, Luis Fernando; Yoshino, Koichiro; Hori, Chiori; Marks, Tim; Polymenakos, Lazaros;
Kummerfeld, Jonathan K.; Galley, Michel; Gao, Xiang

TR2020-029 March 18, 2020

Abstract

This paper provides detailed information about the seventh Dialog System Technology Challenge (DSTC7) and its three tracks aimed to explore the problem of building robust and accurate end-to-end dialog systems. In more detail, DSTC7 focuses on developing and exploring end-to-end technologies for the following three pragmatic challenges: (1) sentence selection for multiple domains, (2) generation of informational responses grounded in external knowledge, and (3) audio visual scene-aware dialog to allow conversations with users about objects and events around them. This paper summarizes the overall setup and results of DSTC7, including detailed descriptions of the different tracks, provided datasets and annotations, overview of the submitted systems and their final results. For Track 1, LSTM-based models performed best across both datasets, allowing teams to effectively handle task variants where no correct answer was present or when multiple paraphrases were included. For Track 2, RNN-based architectures augmented to incorporate facts by using two types of encoders: a dialog encoder and a fact encoder plus using attention mechanisms and a pointer-generator approach provided the best results. Finally, for Track 3, the best model used Hierarchical Attention mechanisms to combine the text and vision information obtaining a 22% better result than the baseline LSTM system for the human rating score. More than 220 participants were registered and about 40 teams participated in the final challenge. 32 scientific papers reporting the systems submitted to DSTC7, and 3 general technical papers for dialog technologies, were presented during the one-day wrap-up workshop at AAAI-19. During the workshop, we reviewed the state-of-the-art systems, shared novel approaches to the DSTC7 tasks, and discussed the future directions for the challenge (DSTC8).

Computer Speech and Language

© 2020 MERL. This work may not be copied or reproduced in whole or in part for any commercial purpose. Permission to copy in whole or in part without payment of fee is granted for nonprofit educational and research purposes provided that all such whole or partial copies include the following: a notice that such copying is by permission of Mitsubishi Electric Research Laboratories, Inc.; an acknowledgment of the authors and individual contributions to the work; and all applicable portions of the copyright notice. Copying, reproduction, or republishing for any other purpose shall require a license with payment of fee to Mitsubishi Electric Research Laboratories, Inc. All rights reserved.

Overview of the seventh Dialog System Technology Challenge: DSTC7

Luis Fernando D'Haro^a, Koichiro Yoshino^b, Chiori Hori^c, Tim K. Marks^c,
Lazaros Polymenakos^d, Jonathan K. Kummerfeld^e, Michel Galley^f, Xiang
Gao^{f,1}

^a*Speech Technology Group. Center for Information Processing and Telecommunications
(IPTC), ETSI Telecomunicación Universidad Politécnica de Madrid, Ciudad
Universitaria, Av. Complutense, 30, 28040 Madrid, Spain*

^b*Nara Institute of Science and Technology, Ikoma, Nara, 6300192, Japan*

^c*Mitsubishi Electric Research Laboratories (MERL), 201 Broadway, Cambridge, MA,
02139, USA*

^d*Alexa Dialog Science, 101 Main Street, Cambridge, MA, 02142, USA*

^e*University of Michigan, 2260 Hayward Street, Ann Arbor, MI 48109, USA*

^f*Microsoft Research, One Microsoft Way, Redmond, WA, 98052, USA*

Abstract

This paper provides detailed information about the seventh Dialog System Technology Challenge (DSTC7) and its three tracks aimed to explore the problem of building robust and accurate end-to-end dialog systems. In more detail, DSTC7 focuses on developing and exploring end-to-end technologies for the following three pragmatic challenges: (1) sentence selection for multiple domains, (2) generation of informational responses grounded in external knowledge, and (3) audio visual scene-aware dialog to allow conversations with users about objects and events around them.

This paper summarizes the overall setup and results of DSTC7, including detailed descriptions of the different tracks, provided datasets and annotations, overview of the submitted systems and their final results. For Track 1, LSTM-based models performed best across both datasets, allowing teams to effectively handle task variants where no correct answer was present or when multiple paraphrases were included. For Track 2, RNN-based architectures augmented to incorporate facts by using two types of encoders: a dialog encoder and a fact encoder plus using attention mechanisms and a

¹Every author has equal contribution. <http://workshop.colips.org/dstc7>

pointer-generator approach provided the best results. Finally, for Track 3, the best model used Hierarchical Attention mechanisms to combine the text and vision information obtaining a 22% better result than the baseline LSTM system for the human rating score.

More than 220 participants were registered and about 40 teams participated in the final challenge. 32 scientific papers reporting the systems submitted to DSTC7, and 3 general technical papers for dialog technologies, were presented during the one-day wrap-up workshop at AAIL-19. During the workshop, we reviewed the state-of-the-art systems, shared novel approaches to the DSTC7 tasks, and discussed the future directions for the challenge (DSTC8).

DSTC7: Dialog Challenge to build more robust and accurate end-to-end dialog systems.

Track 1, Sentence selection for multiple domains, including variations where there are a large number of candidate options, and where the candidate set has zero, one, or multiple correct options.

Track 2, Beyond Chitchat: Generation of informational responses grounded in external knowledge.

Track 3, Audio visual scene-aware dialog systems to allow dynamic conversations about objects and events around users.

10 Overview of the seventh Dialog System Technology
11 Challenge: DSTC7

12 Luis Fernando D'Haro^a, Koichiro Yoshino^b, Chiori Hori^c, Tim K. Marks^c,
13 Lazaros Polymenakos^d, Jonathan K. Kummerfeld^e, Michel Galley^f, Xiang
14 Gao^{f,1}

15 ^a*Speech Technology Group. Center for Information Processing and Telecommunications*
16 *(IPTC), ETSI Telecomunicación Universidad Politécnica de Madrid, Ciudad*
17 *Universitaria, Av. Complutense, 30, 28040 Madrid, Spain*

18 ^b*Nara Institute of Science and Technology, Ikoma, Nara, 6300192, Japan*

19 ^c*Mitsubishi Electric Research Laboratories (MERL), 201 Broadway, Cambridge, MA,*
20 *02139, USA*

21 ^d*Alexa Dialog Science, 101 Main Street, Cambridge, MA, 02142, USA*

22 ^e*University of Michigan, 2260 Hayward Street, Ann Arbor, MI 48109, USA*

23 ^f*Microsoft Research, One Microsoft Way, Redmond, WA, 98052, USA*

24 **Abstract**

25 This paper provides detailed information about the seventh Dialog System
26 Technology Challenge (DSTC7) and its three tracks aimed to explore the
27 problem of building robust and accurate end-to-end dialog systems. In more
28 detail, DSTC7 focuses on developing and exploring end-to-end technologies
29 for the following three pragmatic challenges: (1) sentence selection for multi-
30 ple domains, (2) generation of informational responses grounded in external
31 knowledge, and (3) audio visual scene-aware dialog to allow conversations
32 with users about objects and events around them.

33 This paper summarizes the overall setup and results of DSTC7, including
34 detailed descriptions of the different tracks, provided datasets and annota-
35 tions, overview of the submitted systems and their final results. For Track
36 1, LSTM-based models performed best across both datasets, allowing teams
37 to effectively handle task variants where no correct answer was present or
38 when multiple paraphrases were included. For Track 2, RNN-based archi-
39 tectures augmented to incorporate facts by using two types of encoders: a
40 dialog encoder and a fact encoder plus using attention mechanisms and a
41 pointer-generator approach provided the best results. Finally, for Track 3,
42 the best model used Hierarchical Attention mechanisms to combine the text

¹Every author has equal contribution. <http://workshop.colips.org/dstc7>
Preprint submitted to Computer Speech and Language July 30, 2019

43 and vision information obtaining a 22% better result than the baseline LSTM
44 system for the human rating score.

45 More than 220 participants were registered and about 40 teams partici-
46 pated in the final challenge. 32 scientific papers reporting the systems sub-
47 mitted to DSTC7, and 3 general technical papers for dialog technologies,
48 were presented during the one-day wrap-up workshop at AAAI-19. During
49 the workshop, we reviewed the state-of-the-art systems, shared novel ap-
50 proaches to the DSTC7 tasks, and discussed the future directions for the
51 challenge (DSTC8).

52 *Keywords:*

53 Dialog System Technology Challenge, end-to-end dialog systems, Sentence
54 Selection, Natural Language Generation, Audio Visual Scene-Aware Dialog.

55 1. Introduction

56 The ongoing DSTC series started as an initiative to provide a common
57 testbed for the task of Dialog State Tracking; the first edition was organized
58 in 2013 (Williams et al. (2013)) and used human-computer dialogs in the
59 bus timetable domain. Dialog State Tracking Challenges 2 (Henderson et al.
60 (2014a)) and 3 (Henderson et al. (2014b)) followed in 2014, using more com-
61 plicated and dynamic dialog states for restaurant information in different
62 situations, e.g. state tracking for unseen states, and tested with different do-
63 main data. Dialog State Tracking Challenge 4 (Kim et al. (2017)) and Dialog
64 State Tracking Challenge 5 (Kim et al. (2016)) moved to tracking human-
65 human dialogs in mono- and cross-language settings. Then, for DSTC6 in
66 2017, the challenge focused on end-to-end systems with the aim of minimiz-
67 ing effort on human annotation while exploring more complex and diverse
68 tasks related with dialog systems (Hori et al. (2019c)). For this last edition,
69 DSTC7 in 2018, we focused on scaling the capabilities of the systems, explore
70 multimodal approaches and better use of external information.

71 It is clear that, since its first edition in 2013, the challenge has evolved
72 in several ways. First, from modeling human-computer interactions, then to
73 explore human-human interactions, and finally moving toward complex and
74 more robust end-to-end systems. DSTC has also offered pilot tasks on speech
75 act prediction, spoken language understanding, natural language generation,
76 and end-to-end system evaluation, which expanded interest in the challenge
77 for the dialog and AI research communities. Therefore, given the remarkable

78 success of the first five editions, the complexity of the dialog phenomenon
79 and the interest of the research community in the broader variety of dialog
80 related problems, the DSTC rebranded itself as “Dialog System Technology
81 Challenges” since its sixth edition.

82 For the seventh edition, there were five task proposals. These were dis-
83 cussed during the AAIL-19 workshop, with a focus on how applied proposals
84 were, and how they fit within the larger space of problems of interest to
85 the research community. Three critical issues were raised in the discussion.
86 First, despite the enormous success of the generative approaches used in neu-
87 ral conversation models for response generation, retrieval-based approaches
88 are still essential from a practical point of view (Sentence Selection Track).
89 Second, improving generative approaches is important too in order to allow
90 more response variety considering the dialog context, dialog history, other
91 dialog situations, and grounding the responses by means of external knowl-
92 edge (Sentence Generation Track). The final issue was to extend the dialog
93 systems with complementary multimodal information to allow the system to
94 understand better the context, and allowing the fusion with other research
95 areas; visual dialog is one direction in which information in images is used
96 in the dialog (Audio Visual Scene-Aware Dialog Track). Following this dis-
97 cussion, three tasks were selected for the seventh Dialog System Technology
98 Challenge, as described below.

99 For the Sentence Selection track (described in more detail in section 2),
100 the challenge consists of five sub-tasks, in which systems are given a partial
101 conversation, and they must select the correct next utterance from a short or
102 very large set of candidates, including paraphrases as candidates, or indicate
103 that none of the proposed utterances is correct. This is intended to push the
104 utterance classification task towards real-world problems.

105 For the Sentence Generation track (described in detail in section 3),
106 the goal is to generate informative responses that go beyond chitchat, in
107 this case by injecting informational responses that are grounded in external
108 knowledge (e.g., news stories, or background information such as Wikipedia
109 pages). This task is indented to promote research on fully data-driven re-
110 sponse generation—which has so far been mostly limited to chitchat—by
111 combining the benefits of fully end-to-end approaches with more practical
112 purposes (e.g., informing the users rather than just entertaining them).

113 Finally, in the Audio Visual Scene-aware Dialog track (described in detail
114 in section 4), the goal is to generate system responses in a dialog about an
115 input video. Dialog systems need to understand scenes to have conversa-

116 tions with users about the objects and events around them. In this track,
117 multiple research technologies are integrated including: end-to-end dialog
118 technologies, which generate system responses using models trained from di-
119 alog data; visual question answering (VQA) technologies, which answer to
120 questions about images using learned image features; and video description
121 technologies, in which videos are described/narrated using multimodal infor-
122 mation.

123 *1.1. Workshop summary and future DSTC*

124 The workshop for the Dialog System Technology Challenge (DSTC) was
125 held on January 27, 2019 at Honolulu, Hawaii, USA, collocated with the
126 Thirty-Third AAAI Conference on Artificial Intelligence (AAAI-19). More
127 than 220 participants were registered in one or several of the proposed three
128 tasks; finally, about 40 teams submitted their final results and 32 scientific
129 papers were presented during the workshop, together with 3 general technical
130 papers about dialog systems. We had about 80 pre-registrations for the
131 workshop and more participants joined on-site. The workshop also had many
132 supporting organizations including three sponsors, and an invited talk about
133 Massively Multilingual Dialog and Q&A by Dr. Holger Schwenk.

134 In addition, as part of our efforts to promote the research in dialog tech-
135 nologies, we presented the challenge, tracks, provided data and results during
136 the 2nd NeurIPS workshop on Conversational AI: Today’s Practice and To-
137 morrow’s Potential².

138 Finally, to initiate DSTC8, from November 22, 2018 until January 11,
139 2019 we received up to 7 track proposals for DSTC8³. During the AAAI-
140 19 workshop these proposals were presented to the attendees and then we
141 passed them a survey to know their interest and willingness to participate
142 on each; after the workshop, the following tracks were selected: a) End-to-
143 end Task Completion b) Predicting Responses, c) Audio Visual Scene-Aware
144 Dialog, and d) Schema-Guided State Tracking. This way, we will continue
145 focusing on end-to-end dialog tasks and their application to Dialog Systems
146 in a pragmatic way.

²<http://alborz-geramifard.com/workshops/nips18-Conversational-AI/Main.html>

³For detailed information about each proposal and the selection criteria check:
http://workshop.colips.org/dstc7/dstc8_proposals.html

147 2. Sentence Selection Track

148 Automatic dialogue systems have great potential as a new form of user
149 interface between people and computers. Unfortunately, there are relatively
150 few large resources of human-human dialogues (Serban et al., 2018), which
151 are crucial for the development of robust statistical models. Evaluation also
152 poses a challenge, as the output of an end-to-end dialogue system could
153 be entirely reasonable, but not match the reference, either because it is a
154 paraphrase, or it takes the conversation in a different, but still coherent,
155 direction.

156 In this track, we introduced two new datasets and explored variations in
157 task structure for research on goal-oriented dialogue. One of our datasets was
158 carefully constructed with real people acting in a university student advising
159 scenario. The other dataset was formed by applying a new disentanglement
160 method (Kummerfeld et al., 2018) to extract conversations from an IRC
161 channel of technical help for the Ubuntu operating system. We structured the
162 dialogue problem as next utterance selection, in which participants receive
163 partial dialogues and must select the next utterance from a set of options.
164 Going beyond prior work, we considered larger sets of options, and variations
165 with either additional incorrect options, paraphrases of the correct option,
166 or no correct option at all. These changes push the next utterance selection
167 task towards real-world dialogue.

168 This task is not a continuation of prior DSTC tasks, but it is related to
169 tasks 1 and 2 from DSTC6 (Perez et al., 2017; Hori and Hori, 2017a). Like
170 DSTC6 task 1, our task considers goal-oriented dialogue and next utterance
171 selection, but our data is from human-human conversations, whereas theirs
172 was simulated. Like DSTC6 task 2, we use online resources to build a large
173 collection of dialogues, but their dialogues were shorter (2 - 2.5 utterances
174 per conversation) and came from a more diverse set of sources (1,242 twitter
175 customer service accounts, and a range of films).

176 Below we provide an overview of (1) the task structure, (2) the datasets,
177 (3) the evaluation metrics, and (4) system results. Twenty teams partici-
178 pated, with one clear winner, scoring the highest on all but one sub-task.
179 The data and other resources associated with the task have been released⁴ to
180 enable future work on this topic and to make accurate comparisons possible.

⁴<https://ibm.github.io/dstc7-noesis/public/index.html>

181 *2.1. Task*

182 This task pushed the state-of-the-art in goal-oriented dialogue systems in
183 four directions deemed necessary for practical automated agents, using two
184 new datasets. We sidestepped the challenge of evaluating generated utter-
185 ances by formulating the problem as next utterance selection, as proposed
186 by Lowe et al. (2015). At test time, participants were provided with partial
187 conversations, each paired with a set of utterances that could be the next
188 utterance in the conversation. Systems needed to rank these options, with
189 the goal of placing the true utterance first. Prior work used sets of 2 or 10
190 utterances. We make the task harder by expanding the size of the sets, and
191 considered several advanced variations:

192 **Subtask 1** 100 candidates, including 1 correct option.

193 **Subtask 2** 120,000 candidates, including 1 correct option (Ubuntu data
194 only).

195 **Subtask 3** 100 candidates, including 1-5 correct options that are paraphrases
196 (Advising data only).

197 **Subtask 4** 100 candidates, including 0-1 correct options.

198 **Subtask 5** The same as subtask 1, but with access to external information.

199 These subtasks push the capabilities of systems. In particular, when
200 the number of candidates is small (2-10) and diverse, it is possible that
201 systems are learning to differentiate topics rather than learning dialogue. Our
202 variations move towards a task that is more representative of the challenges
203 involved in dialogue modeling.

204 As part of the challenge, we provided a baseline system that implemented
205 the Dual-Encoder model from Lowe et al. (2015). This lowered the barrier
206 to entry, encouraging broader participation in the task.

207 *2.2. Data*

208 We used two datasets containing goal-oriented dialogues between two
209 participants, but from very different domains. This challenge introduced the
210 two datasets, and we kept the test set answers secret until after the challenge.⁵

⁵The entire datasets are now publicly available at <https://ibm.github.io/dstc-noesis/public/index.html>

```

210 10:30 <elmaya> is there a way to setup grub to not press the esc button
                for the menu choices?
211 10:31 <scaroo> elmaya, edit /boot/grub/ menu.lst and comment the
                "hidemenu" line
212 10:32 <scaroo> elmaya, then run grub -install
213 10:32 <scaroo> grub-install
214 10:32 <elmaya> thanls scaroo
215 10:32 <elmaya> thanks

```

Figure 1: Example Ubuntu dialogue before our pre-processing.

211 To construct the partial conversations we randomly split each conversation.
 212 Incorrect candidate utterances are selected by randomly sampling utterances
 213 from the rest of the dataset. For subtask 3 (paraphrases), the incorrect
 214 candidates are sampled with paraphrases as well. For subtask 4 (no correct
 215 option sometimes), twenty percent of examples were randomly sampled and
 216 the correct utterance was replaced with an additional incorrect one.

217 Along with the datasets we provided additional sources of information
 218 that were specific to each dataset. Participants were able to use the provided
 219 knowledge sources as is, or automatically transform them to appropriate
 220 representations (e.g. knowledge graphs, continuous embeddings, etc.) that
 221 were integrated with end-to-end dialogue systems so as to increase response
 222 accuracy.

223 2.2.1. *Ubuntu*

224 We constructed one dataset from the Ubuntu Internet Relay Chat (IRC)
 225 support channel, in which users help each other to resolve technical problems
 226 related to the Ubuntu operating system. We consider only conversations in
 227 which one user asks a question and another helps them resolve their problem.
 228 We extracted conversations from the channel using the conversational dis-
 229 entanglement method described by Kummerfeld et al. (2018), trained with
 230 manually annotated data using Slate (Kummerfeld, 2019).^{6,7} See Kummer-
 231 feld et al. (2018) for detailed analysis of the extraction process. At a high

⁶Previously, Lowe et al. (2015) extracted conversations from the same IRC logs, but with a heuristic method. Kummerfeld et al. (2018) showed that the heuristic was far less effective than a trained statistical model.

⁷The specific model used in DSTC 7 track 1 is from an earlier version of Kummerfeld et al. (2018), as described in the ArXiv preprint and released as the C++ version.

232 level, we used a feedforward neural network that considers each message in
233 the logs and predicts which earlier message it is a response to. This forms a
234 structure in which each connected component is a single conversation. The
235 manual annotation of the data had a convention that when a user asks a
236 question that starts a new conversation, which makes it clear who is asking
237 for help and who is providing it.

238 We further applied several filters to increase the quality of the extracted
239 dialogues: (1) the first message must not be directed, (2) there are exactly
240 two participants (a questioner and a helper), not counting the channel bot,
241 (3) no more than 80% of the messages are by a single participant, and (4)
242 there are at least three turns. This approach produced 135,000 conversations,
243 and each was cut off at different points to create the necessary conversations
244 for all the subtasks. In all cases, the cutoff point was chosen to ensure there
245 were at least three prior turns of dialogue.

246 Figure 1 shows an example dialogue from the dataset. For the actual
247 challenge we identify the users as ‘speaker_1’ (the person asking the question)
248 and ‘speaker_2’ (the person answering), and removed usernames from the
249 messages (such as ‘elmaya’ in the example). We also combined consecutive
250 messages from a single user, and always cut conversations off so that the
251 last speaker was the person asking the question. This meant systems were
252 learning to behave like the helpers, which fits the goal of developing a dialogue
253 system to provide help.

254 For subtask 5, additional data was provided in the form of manual pages.
255 These provide information on commands that are frequently mentioned in
256 the Ubuntu technical support conversations.

257 *2.2.2. Advising*

258 Our second dataset is based on an entirely new collection of dialogues in
259 which university students are being advised which classes to take. These were
260 collected at the University of Michigan with IRB approval. Pairs of Michigan
261 students play-acted the roles of a student and an advisor. We provided
262 a persona for the student, describing the classes they had taken already,
263 what year of their degree they were in, and several types of class preferences
264 (workloads, class sizes, topic areas, time of day, etc.). Advisors did not
265 know the student’s preferences, but did know what classes they had taken,
266 what classes were available, and which were suggested (based on aggregate
267 statistics from real student records). The data was collected over a year,
268 with some data collected as part of courses in NLP and social computing,

Student Hi professor, I am looking for courses to take. Do you have any suggestions?

Advisor What topic do you prefer, computer science or electrical engineering?

Student I prefer electrical engineering.

Advisor Based on your background, I would like to suggest you take one of the two courses: EECS 550 Information Theory and EECS 551: Matrix Methods for Signal Processing, Data Analysis and Machine Learning FA 2012

Student Can you describe a little bit about EECS 550?

Advisor This course contains a lot of concepts about source, channel, rate of transformation of information, etc.

Student Sounds interesting. Do you know the class size of this course?

Advisor This is a relatively small class and the average size of it is around 12.

Student I would prefer class with larger class size. What is EECS 551 about?

Advisor This course is about theory and application of matrix methods to signal processing, data analysis and machine learning

Student What is the course size of EECS 551?

Advisor It is around 71

Student I would take EECS 551. Thanks professor!

Advisor You are welcome!

Student Hello, I need some courses. What would you recommend?

Student Hi. I am looking for courses. Could you recommend some?

Student Hi professor, what courses would you suggest for me to take?

Student Professor, what courses do you advise I take?

Student Hi professor, I look for suggestions to take courses. Do you have any?

Student Hi professor, do you have any suggestions to search and take courses?

Advisor Is electrical engineering your main interest, or computer science?

Advisor Do you think you would like computer science or electrical engineering best?

Advisor computer science or electrical engineering, which do you prefer?

Advisor Do u prefer computer science or do u prefer electrical engineering?

Advisor Which subject is more interesting to you, computer science or electrical engineering?

Advisor Do you prefer computer science or electrical engineering?

Figure 2: Example Advising dialogue and paraphrases of the first two utterances.

Property	Advising	Ubuntu
Dialogues	815	135,078
Utterances / Dialogue	18.3	10.0
Tokens / Utterance	9.8	9.9
Utterances / Unique utt.	1.1	1.1
Tokens / Unique tokens	50.8	22.9

Table 1: Comparison of the diversity of the complete underlying datasets (train, dev, test, and unused). Advising is smaller, has longer conversations, and more token diversity. Tokens are based on splitting on whitespace.

269 and some collected with paid participants.

270 In the shared task, we provide all of this information - student pref-
271 erences, and course information - to participants. 815 conversations were
272 collected, and then the data was expanded by collecting 82,094 paraphrases
273 using the crowdsourcing approach described by Jiang et al. (2017). This in-
274 volved asking each worker for multiple paraphrases, with carefully designed
275 examples that guided them towards creative edits that were still correct. Of
276 this data, 500 conversations were used for training, 100 for development, and
277 100 for testing. The remaining 115 conversations were used to create a large
278 pool of utterances. This pool was then used as a source of negative candi-
279 date sentences in the candidate sets. For the test data, 500 conversations
280 were constructed by cutting the conversations off at 5 points and using para-
281 phrases to make 5 distinct conversations. The training data was provided in
282 two forms. First, the 500 training conversations with a list of paraphrases
283 for each utterance, which participants could use in any way. Second, 100,000
284 partial conversations generated by randomly selecting paraphrases for every
285 message in each conversation and selecting a random cutoff point.

286 Two versions of the test data were provided to participants. A mistake
287 led to the first version of the test set drawing from both training and test
288 dialogues, rather than using just the test dialogues. During the challenge this
289 issue was identified and a corrected version was released to all participants.
290 Results on both sets were included in the initial task summary, but we only
291 include the final set here and encourage all future work to only consider the
292 second test set.

293 *2.2.3. Comparison*

294 Table 1 provides statistics about the two raw datasets. The Ubuntu
295 dataset is based on several orders of magnitude more conversations, but they
296 are automatically extracted, which means there are errors (conversations that
297 are missing utterances or contain utterances from other conversations). Both
298 have similar length utterances, but these values are on the original Ubuntu
299 dialogues, before we merge consecutive messages from the same user. The
300 Advising dialogues contain more messages on average, but the Ubuntu dia-
301 logues cover a wider range of lengths (up to 118 messages). Interestingly, the
302 diversity in tokens varies substantially, while utterance lengths and utterance
303 diversity are similar.

304 *2.3. Results*

305 Twenty teams submitted entries for at least one subtask. Additional
306 external resources were not permitted, with the exception of pre-trained em-
307 beddings that were publicly available prior to the release of the data.

308 *2.3.1. Participants*

309 Table 2 presents a summary of approaches teams used. One clear trend
310 was the use of the Enhanced LSTM model (ESIM, Chen et al., 2017), though
311 each team modified it differently as they worked to improve performance on
312 the task. Other approaches covered a wide range of neural model compo-
313 nents: Convolutional Neural Networks, Memory Networks, the Transformer,
314 Attention, and Recurrent Neural Network variants. Two teams used ELMo
315 word representations (Peters et al., 2018), while three constructed ensembles.
316 Several teams also incorporated more classical approaches, such as TF-IDF
317 based ranking, as part of their system.

318 We provided a range of data sources in the task, with the goal of enabling
319 innovation in training methods. Six teams used the external data, while four
320 teams used the raw form of the Advising data. The rules did not state
321 whether the validation data could be used as additional training data at test
322 time, and so we asked each team what they used. As Table 2 shows, only
323 four teams trained their systems with the validation data.

324 *2.3.2. Metrics*

325 We considered a range of metrics when comparing models. Following
326 Lowe et al. (2015), we use Recall@N, where we count how often the correct
327 answer is within the top N specified by a system. In prior work, there were

328 either 2 or 10 candidates (including the correct one), and N was set at 1, 2,
329 or 5. Our sets are larger, with 100 candidates, and so we considered larger
330 values of N: 1, 10, and 50. 10 and 50 were chosen to correspond to 1 and 5 in
331 prior work (the expanded candidate set means they correspond to the same
332 fraction of the space of options). We also considered a widely used metric
333 from the ranking literature: Mean Reciprocal Rank (MRR). For subtask 3
334 we measured Mean Average Precision (MAP) since there are multiple correct
335 utterances in the set. Finally, for subtask 4, participants had to return 101
336 values, the extra one being the value ‘NONE’, to indicate that no valid answer
337 was present.

338 To determine a single winner for each subtask, we used the mean of
339 Recall@10 and MRR, as presented in Table 3.

340 *2.3.3. Discussion*

341 Table 3 presents the overall scores for each team on each subtask, ordered
342 by teams’ average rank. Team 3 consistently scored highest, winning all but
343 one subtask. For details of their approach, see Chen and Wang (2019).
344 Looking at individual metrics, they had the best score 75% of the time on
345 Ubuntu and all of the time on the final Advising test set. The subtask they
346 were beaten on was Ubuntu-2, in which the set of candidates was drastically
347 expanded. Team 10 did best on that task, indicating that their extra filtering
348 step provided a key advantage. They filtered the 120,000 sentence set down
349 to 100 options using a TF-IDF based method, then applied their standard
350 approach to that set. For details of the method, see Ganhotra et al. (2019).

351 *Subtasks.*

- 352 1. The first subtask drew the most interest, with every team participating
353 in it for one of the datasets. Performance varied substantially, covering
354 a wide range for both datasets, particularly on Ubuntu.
- 355 2. As expected, subtask 2 was more difficult than task 1, with consistently
356 lower results. However, while the number of candidates was increased
357 from 100 to 120,000, performance reached as high as half the level of
358 task 1, which suggests systems could handle the large set effectively.
- 359 3. Also as expected, results on subtask 3 were slightly higher than on
360 subtask 1. Comparing MRR and MAP it is interesting to see that
361 while the ranking of systems is the same, in some cases MAP was
362 higher than MRR and in others it was lower.

- 363 4. For both datasets, results on subtask 4, where the correct answer was
364 to choose no option 20% of the time, are generally similar. On average,
365 no metric shifted by more than 0.016, and some went up while others
366 went down. This suggests that teams were able to effectively handle
367 the added challenge.
- 368 5. Finally, on subtask 5 we see some slight gains in performance, but
369 mostly similar results, indicating that effectively using external re-
370 sources remains a challenge.

371 *Advising Test Sets.* We compared results on the two versions of the test set
372 (one which had overlap with the source dialogues from training, and the
373 other with entirely distinct dialogues). Removing overlap made the task
374 considerably harder, though more realistic. In general, system rankings were
375 not substantially impacted, with the exception of team 17, which did better
376 on the original dataset. This may relate to their use of a memory network
377 over the raw advising data, which may have led the model to match test
378 dialogues with their corresponding training dialogues.

379 *Metrics.* Finally, we compared the metrics. In 39% of cases a team’s ranking
380 is identical across all metrics, and in 34% there is a difference of only one
381 place. The maximum difference is 5, which occurred once, between team 6’s
382 results in the final Advising results, where their Recall@1 result was 8th, their
383 Recall@10 result was 11th and their Recall@50 result was 13th. Comparing
384 MRR and Recall@N, the MRR rank is outside the range of ranks given by the
385 recall measures 9% of the time (on Ubuntu and the final Advising evaluation).

386 2.4. Future Work

387 This task provides the basis for a range of interesting new directions.
388 We randomly selected negative options, but other strategies could raise the
389 difficulty, for example by selecting very similar candidates according to a
390 simple model. For evaluation, it would be interesting to explore human
391 judgements, since by expanding the candidate sets we are introducing options
392 that are potentially reasonable.

393 This work has been extended in several direction by a follow-up task at
394 DSTC 8. In particular, the setting was expanded to include conversations
395 with more than two participants. One subtask also explores the challenge
396 of selecting responses in the raw channel, where multiple conversations are
397 occurring at once. These pose additional challenges and bring the setting

398 closer to the real world. The data has also been improved, by using an im-
399 proved version of the disentanglement algorithm that extracts higher quality
400 conversations.

401 2.5. Conclusion

402 This task introduced two new datasets and three new variants of the next
403 utterance selection task. Twenty teams attempted the challenge, with one
404 clear winner. The datasets are being publicly released, along with a baseline
405 approach, in order to facilitate further work on this task. This resource will
406 support the development of novel dialogue systems, pushing research towards
407 more realistic and challenging settings.

408 3. Sentence Generation Track

409 Recent work (Ritter et al., 2011; Sordoni et al., 2015; Shang et al., 2015;
410 Vinyals and Le, 2015; Serban et al., 2016, etc.) has shown that conversa-
411 tional models can be trained in a completely end-to-end and data-driven
412 fashion, without any hand-coding. However, prior work has mostly focused
413 to chitchat, as that is a common feature of messages in the social media data
414 (e.g., Twitter (Ritter et al., 2011)) used to train these systems. Such end-to-
415 end neural conversation systems have a tendency to produce responses that
416 are conversationally appropriate, but that are also often bland (Li et al.,
417 2016a; Gao et al., 2019), purely chatty, and lacking entities and factual con-
418 tent. On the other end, goal-oriented dialog systems have the ability to
419 inject entities and facts into responses, but often at the cost of significant
420 hand-coding (e.g., slot filling) and this hand-crafting is often specific to the
421 domain or task. We argue that dialog shouldn’t necessarily be either com-
422 pletely goal-oriented or completely chitchat. This is often reflected in real
423 human-human data, which often combines the two genres.

424 To effectively move beyond chitchat and produce system responses that
425 are both substantive and “useful”, fully data-driven models need grounding
426 in the real world and access to external knowledge (textual or structured). To
427 do so, the Sentence Generation task was inspired by the *knowledge-grounded*
428 conversational framework of (Ghazvininejad et al., 2018; Qin et al., 2019),
429 which combines conversational input and textual data from the user’s envi-
430 ronment (here, a web page that is discussed). Such a framework maintains
431 the benefit of fully data-driven conversation while attempting to get closer

432 to task-oriented scenarios, with the goal of informing and helping the users
433 and not just entertaining them.

434 3.1. Task definition

435 The task follows the data-driven framework established in 2011 by Rit-
436 ter et al. (2011), which avoids hand-coding any linguistic, domain, or task-
437 specific information (e.g., there are no explicit dialog act or slots). In the
438 knowledge-grounded setting of (Ghazvininejad et al., 2018; Qin et al., 2019),
439 that framework is extended as each system input consists of two parts:

- 440 • **Conversational input:** Similar to DSTC6 Track 2 (Hori and Hori,
441 2017b), all preceding turns of the conversation are available to the
442 system. For practical purposes, we truncate the context to the K most
443 recent turns.
- 444 • **Contextually-relevant “facts”:** The system is given text that is
445 relevant to the context of the conversation, in this case a web page.
446 This text is distinct from conversational data, and is extracted from
447 external knowledge sources such as Wikipedia or news web sites.

448 From this input, the task is to produce a response that is (1) conversa-
449 tionally appropriate and relevant, as well as (2) informative and interesting.
450 The evaluation setup is presented in Section 3.4, which includes a human
451 evaluation of these two qualities (“Relevance” and “Interest”, respectively).

452 3.2. Data

453 We extracted conversation threads from Reddit data, which is particularly
454 well suited for grounded conversation modeling. Indeed, Reddit conversations
455 are organized around submissions, where each conversation is typically initi-
456 ated with a URL to a web page (grounding) that defines the subject of the
457 conversation. An example of the data is shown in Table 4. For this task, we
458 restrict ourselves to submissions that contain exactly one URL and a title. To
459 reduce spamming and offensive language and improve the overall quality of
460 the data, we restricted our grounded dataset to 226 web domains and to 178
461 high-quality Reddit topics (i.e., “subreddits”). We also imposed constraints
462 on turn length similar to those in place in Twitter (e.g., responses must be
463 less than 280 characters), in order to ensure that dialogue turns are con-
464 versational and not long monologues. This filtering yielded about 3 million

465 conversational responses and 20 million facts.⁸ We split the data into train,
466 validation and test, with the following month ranges for these different sets:
467 years 2011-2016 for train, Jan-Mar 2017 for validation, and the rest of 2017
468 for test. For the test set, we selected conversational turns for which 6 or more
469 responses were available, in order to create a multi-reference test set. Given
470 other filtering criteria such as turn length, this yielded a 5-reference test set
471 of size 2208 (For each instance, we set aside one of the 6 human responses to
472 assess human performance on this task). More information about the data
473 can be found in Qin et al. (2019), which introduced this dataset. All code
474 and data can also be found on the DSTC Track 2 page,⁹ which makes data
475 extraction, baseline, and evaluation code available, and lets anyone recreate
476 the training, development, validation and test sets.

477 3.3. Submitted Systems

478 The submitted systems include sequence-to-sequence models (Sordoni
479 et al., 2015; Shang et al., 2015; Vinyals and Le, 2015) with memory network
480 and related models (Weston et al., 2015; Sukhbaatar et al., 2015), copy-based
481 mechanism (See et al., 2017; Gu et al., 2016; He et al., 2017), hierarchical
482 model (Serban et al., 2016), attention mechanism (Bahdanau et al., 2015),
483 and variational model (Kingma and Welling, 2013). The following is a brief
484 summary of the systems based on system descriptions and private commu-
485 nication:

- 486 • **TeamA:** Details of this systems are unknown to us as a system de-
487 scription was not submitted.
- 488 • **TeamB:** It is a sequence-to-sequence model with a copying mechanism
489 (See et al., 2017) from both the conversation history and facts. A
490 modified beam search with some semantic clustering is proposed to
491 discourage bland or meaningless responses.
- 492 • **TeamC:** It is a sequence-to-sequence modeling the skeleton of dialog
493 response for pretraining, then fine-tuned with a Memory Network en-

⁸We could have easily increased the number of web domains to create a bigger dataset, but we aimed to make the task relatively accessible for participants with limited computing resources.

⁹<https://github.com/mgalley/DSTC7-End-to-End-Conversation-Modeling>

494 coder (Sukhbaatar et al. (2015)) that utilizes retrieved top-10 related
495 facts.

496 • **TeamD:** This system consists of a Memory-augmented Hierarchical
497 Encoder-Decoder (MHRED) that extends (Serban et al., 2016), a sen-
498 tence selection module to retrieve facts, and a reranker.

499 • **TeamF:** It is a variational generative model with a joint attention
500 mechanism conditioning on the contexts and textual facts.

501 • **TeamG:** It is a variational generative model. Contexts (and response
502 at the training stage) are encoded to extract textual fact information
503 using an attention mechanism.

504 3.4. Evaluation

505 We evaluated response quality using both automatic and human evalu-
506 ation. Since we are not considering task-oriented dialog, there is no pre-
507 specified task and therefore no extrinsic way of measuring task success. In-
508 stead, we performed a per-response human evaluation judging each system
509 response using crowdsourcing:

510 • **Relevance:** This evaluation criterion measures whether the system
511 response is conversationally appropriate and relevant to the given K
512 immediately preceding turns (to reduce the judges’ cognitive load we
513 set K as 2). Grounding in external sources is not involved in this judge.

514 • **Interest:** This evaluation criterion asks whether the produced response
515 is interesting and informative given the document provided by the URL.
516 To reduce cognitive load, we only considered URLs with named anchors
517 (i.e., prefixed with ‘#’ in the URL) and only a snippet of the document
518 immediately following that anchor is provided to the crowdworkers.
519 Note that models could use full web pages as input.

520 Both evaluation criteria were scored on a 5-point Likert scale, and finally
521 combined the two judgments with equal weights.

522 In order to provide participants with preliminary results to include in their
523 system descriptions, we also performed automatic evaluation using standard
524 machine translation metrics, including BLEU (Papineni et al., 2002), ME-
525 TEOR (Lavie and Agarwal, 2007), and NIST (Doddington, 2002). NIST is a

526 variant of BLEU that weights n -gram matches by their information gain, i.e.,
527 it indirectly penalizes uninformative n -grams such as “I don’t” and “don’t
528 know”. The final ranking of the systems was based only on human evaluation
529 scores.

530 3.5. Results

531 3.5.1. Automatic Evaluation

532 The Generation Task received 26 system submissions from 7 teams. In
533 addition to these systems, we also evaluated a “human” system (one of the
534 six human references set aside for evaluation) and three baselines: a seq2seq
535 baseline, a “random_human” baseline (which randomly selects human re-
536 sponses from the training data), and a constant baseline (which always re-
537 sponds “I don’t know what you mean.”).¹⁰ The reason for including a con-
538 stant baseline is that such a deflective response generation system can be
539 surprisingly competitive, at least when evaluated on automatic metrics (e.g.,
540 BLEU). While the idea of such a constant baseline is relatively new, it is
541 inspired by the idea that open-domain conversational systems trained end-
542 to-end have a tendency to produce outputs that are relatively constant (Li
543 et al., 2016b), such as “I don’t know.” The main automatic score results are
544 shown in Table 5, and the findings for each of the metrics are as follows:

- 545 • **BLEU-4:** When evaluated on 5 references, the constant baseline,
546 which always responds deflectively, does surprisingly well (2.87%) and
547 outperforms all the submitted systems (ranging from 1.01% to 1.83%),
548 and is only outperformed by humans. In further analysis, we found
549 that reducing the number of references to one solved the problem, as
550 almost all the systems were able to outperform the baseline accord-
551 ing to single-reference BLEU. We suspect this deficiency of BLEU *with*
552 *many references*, previously noted in Vedantam et al. (2015a), to be
553 due to its parameterization as a precision metric. For example, if one
554 of the gold responses happens to be “I don’t know what you mean”,
555 the constant baseline gets a maximum score for that instance, irrespec-
556 tively of all other references. Thus, this biases the metric towards very
557 bland responses, as often at least one of the 5 references is somewhat
558 deflective (e.g., contains “I don’t know”). Based on these observations,

¹⁰This constant response was greedily selected to optimize a combination of BLEU, NIST, and METEOR on a held-out set.

559 we recommend to use single-reference BLEU instead of multi-reference
560 BLEU for future DSTC tasks similar to this task, as the former gave
561 much more meaningful results.

562 • **NIST-4:** The NIST score weights n -gram matches by their informa-
563 tion gain, and effectively penalizes common n -grams such as “I don’t
564 know”, which alleviates the problem with multi-reference BLEU men-
565 tioned above. None of the baselines is competitive with the top systems
566 according to NIST-4, even when using 5 references. This suggests that
567 NIST might be a more suitable metric than BLEU when dealing with
568 multi-reference test sets, and it penalizes bland responses. Note that
569 the “Random.Human” system does relatively well according to NIST-4,
570 but this is probably due to the fact that this random baseline selects *hu-*
571 *man* sentences randomly from the training data, and human responses
572 generally contain n -grams with more information content than machine
573 generated n -grams.

574 • **METEOR:** This metric suffers from the same problem as BLEU-4,
575 as the constant baseline performs very well on that metric and outper-
576 forms all submitted primary systems but one. We suspect this is due
577 to the fact that METEOR (as BLEU) does not consider information
578 gain in its scoring.

579 Table 5 also provides unigram and bigram diversity scores as defined in Li
580 et al. (2016c), which are important to qualify the performance of some of the
581 systems and baselines. Indeed, a high BLEU score (e.g., constant baseline)
582 can be a consequence of very bland and uninformative output.

583 In future work, we will also consider comparing these metrics against
584 CIDEr (Vedantam et al. (2015b)), AM-FM (D’Haro et al. (2019), Banchs
585 et al. (2015)) Embedding Average cosine similarity, Skip-Thoughts cosine
586 similarity, and other metrics used before in dialogue (Sharma et al. (2017)).

587 3.5.2. Human Evaluation

588 We limited evaluation to a sample of 1000 conversations and only used
589 primary systems due to the cost of crowd-sourcing. All systems were evalu-
590 ated with the same set of conversations, and results are displayed in Table 6.

591 Each output was judged by 3 randomly-assigned judges for Relevance and

592 Interest using a 5-point Likert scale. After removing spamming,¹¹ inter-rater
593 agreement on a converted 3-way scale was fair, as indicated by Fleiss’ Kappa
594 at 0.39 for Relevance and 0.38 for Interest. As expected, the constant baseline
595 performed moderately well on Relevance (2.60), but received a relatively low
596 Interest score (constant: 2.32). The best system returned a composite score
597 of 2.93 (Relevance: 2.99, Interest: 2.87), but is still below the human level
598 of 3.55 (Relevance: 3.61, Interest: 3.49).

599 Finally, we assess the level of correlation between automatic and hu-
600 man scores for this task, to help determine whether it would be appropri-
601 ate to rely mostly on automatic evaluation in future end-to-end response
602 generation tasks similar to DSTC Track 2. We computed system-level cor-
603 relation between overall human scores (i.e., relevance+interest) on the one
604 hand, and each of the individual main metric on the other hand (i.e., ei-
605 ther BLEU-4, NIST-4, and METEOR).¹² We found that automatic metrics’
606 Spearman rank correlation coefficients (ρ) computed against human scores
607 to be quite promising, with $\rho = 0.535$ for BLEU-4, $\rho = 0.650$ for METEOR,
608 and $\rho = 0.669$ for NIST-4. As Table 5 suggests that BLEU-4 and NIST-4
609 tend to complement each other (with NIST-4 giving high scores to diverse
610 responses, and BLEU-4 penalizing them), we also computed the correlation
611 between the unweighted linear combination of these 3 metrics on one hand
612 (Figure 3), and overall human scores on the other hand: this yield Spear-
613 man’s $\rho = 0.754$. While this result indicates a rather strong correlation
614 between human ratings and automatic metrics for this task, it is probably
615 not strong enough to warrant bypassing human evaluation altogether, espe-
616 cially given the small sample size of this correlation analysis. Nonetheless,
617 we consider this result to be relatively positive, as we believe it would pro-
618 vide participants of future end-to-end responses generation tasks a quick and
619 relatively decent substitute to human judgment in their day-to-day (i.e., not

¹¹We removed annotation of judges suspected to be spammers if their rating diverged significantly from the mean ratings of the other judges (i.e., correlation coefficient close to zero.) Such a situation is usually a sign that the judge is either rating deterministically without looking at the task (e.g., always selecting the first option in the list or ratings) or is rating randomly.

¹²Note that we computed system-level rather than sentence-level correlation, as the BLEU-4 and NIST-4 metrics were designed to be computed at a corpus rather than sentence level, as some of their underlying statistics (e.g., 4-gram matches) cannot be reliably computed on single turns or sentences.

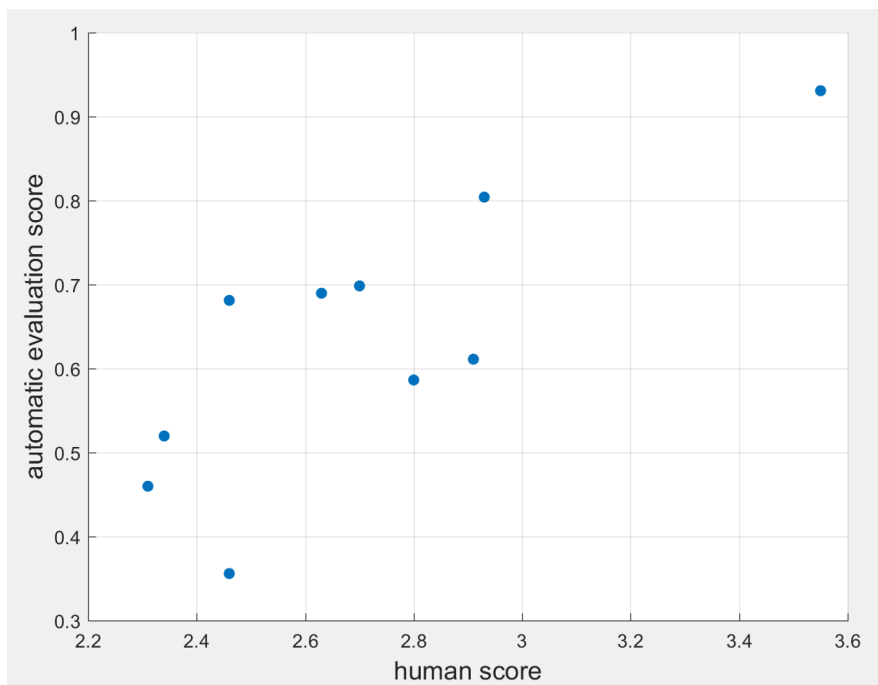


Figure 3: System-level correlation between overall human scores (relevance+interest) and automatic evaluation (unweighted linear combination of BLEU-4, NIST-4, and METEOR).

620 final) system performance evaluations.

621 3.6. Summary

622 The sentence generation task challenged participants to produce interest-
 623 ing and informative end-to-end conversational responses that drew on tex-
 624 tual background knowledge. In this respect, the task was significantly more
 625 challenging than the DSTC6 task that was focused on the conversational di-
 626 mensions of response generation. In general, competing system outputs were
 627 judged by humans to be more relevant and interesting than our constant and
 628 random baselines. It is also clear, however, that the quality gap between
 629 human and system responses is substantial, indicating that there is consid-
 630 erable space for research in future algorithmic improvements. For the future
 631 work, one line of investigation will be to explore the effect of other mech-
 632 anism to extract information from the textual grounding, such as off-the-shelf
 633 machine reading models including BERT Devlin et al. (2018). Multimodal

634 grounding is another line of future work.

635 4. Audio Visual Scene-aware Dialog Track

636 In this track, we consider a new research target: a dialog system that can
637 discuss dynamic scenes with humans. This lies at the intersection of research
638 in natural language processing, computer vision, and audio processing. As
639 described above, end-to-end dialog modeling using paired input and output
640 sentences has been proposed as a way to reduce the cost of data prepara-
641 tion and system development. Such end-to-end approaches have been shown
642 to better handle flexible conversations by enabling model training on large
643 conversational datasets (Vinyals and Le, 2015; Hori et al., 2019c). However,
644 current dialog systems cannot understand a scene and have a conversation
645 about what is going on in it. To develop systems that can carry on a con-
646 versation about objects and events taking place around the machines or the
647 users, systems need to understand not only the dialog history but also the
648 video and audio information in the scene. In the field of computer vision,
649 interaction with humans about visual information has been explored in *visual*
650 *question answering* (VQA) by Antol et al. (2015) and *Visual Dialog* by Das
651 et al. (2017). These tasks have been the focus of intense research, aiming to
652 (1) generate answers to questions about things and events in a single static
653 image and (2) hold a meaningful dialog with humans about an image using
654 natural, conversational language in an end-to-end framework. While VQA
655 and visual dialog take significant steps towards human-machine interaction,
656 they only consider a single static image. Most real-world scenarios, such as
657 helping visually impaired users or intelligent home assistants, involve time-
658 varying information. Thus, they need to be able to process video information
659 to understanding the content and temporal dynamics of a scene. To capture
660 the semantics of dynamic scenes, recent research has focused on *video de-*
661 *scription*. The state of the art in video description uses multimodal fusion
662 to combine different input modalities (feature types), such as the attention-
663 based fusion of spatio-temporal motion features and audio features proposed
664 by Hori et al. (2017).

665 Since the recent revolution of neural network models allows us to combine
666 different modules into a single end-to-end differentiable network, this frame-
667 work allow us to build scene-aware dialog systems by combining end-to-end
668 dialog and multimodal video description approaches. We can simultaneously

669 input video features and user utterances into an encoder-decoder-based sys-
670 tem whose outputs are natural-language responses.

671 To advance this goal, we introduce a new dataset of human dialogues
672 about videos. As the subject matter of Audio Visual Scene-aware Dialog
673 (AVSD), we used the short video clips of the Charades dataset (Sigurdsson
674 et al., 2016): simple videos of real people performing everyday actions in
675 real-world settings, with natural audio. The baseline system we provided in-
676 corporated technologies for video description into an end-to-end dialog sys-
677 tem (Hori et al., 2018a). We made the dataset, code, and model publicly
678 available for a new Audio Visual Scene-Aware Dialog (AVSD) Challenge at
679 DSTC7.

680 4.1. Task definition

681 In this track, the system must generate responses to a user input in the
682 context of a given dialog. The target of *VQA* and *Visual Dialog* is sentence
683 selection based on information retrieval. For real-world application, however,
684 spoken dialog systems cannot simply select from a small set of pre-determined
685 sentences. Instead, they need to immediately output a response to a user
686 input. For this reason, in this track we focus on sentence generation rather
687 than sentence selection. In this track, the system’s task is to use a dialog
688 history (the previous rounds of questions and answers in a dialog between
689 user and system) and (optionally) a brief video script, plus (in one version of
690 the task) the visual and audio information from the input video, to answer a
691 next question about the video. There are two tasks, each with two versions
692 (a and b):

693 **Task 1: Video and Text** (a) Using the video and text training data pro-
694 vided but no external data sources, other than publicly available pre-
695 trained feature extraction models (b) Also using external data for train-
696 ing.

697 **Task 2: Text Only** (a) Do not use the input videos nor their audio tracks
698 for training or testing. Use only the text training data (dialog history
699 and video script) provided. (b) Any publicly available text data may
700 be used for training.

701 4.2. Data

702 To set up the Audio Visual Scene-Aware Dialog (AVSD) track, we col-
703 lected (in Alamri et al. (2018a)) text-based dialogs about short videos from

704 the Charades dataset (Sigurdsson et al., 2016)¹³, which consists of untrimmed
705 and multi-action videos along with a brief script for each video. The data
706 collection paradigm for dialogs was similar to the one described by Das et al.
707 (2016), in which for each image, two parties interacted via a text interface
708 to yield a dialog. In Das et al. (2016), each dialog consisted of a sequence
709 of questions and answers about an image. In our audio visual scene-aware
710 dialog case, two parties had a discussion about events in a video. One of
711 the two parties played the role of an answerer who had already watched the
712 video and read the video script. The answerer answered questions asked by
713 their counterpart, the questioner. The questioner was not allowed to watch
714 the video but was able to see the first, middle, and last frames of the video
715 as single static images. The two had 10 rounds of Q and A, in which the
716 questioner asked about the events that happened in the video. At the end,
717 the questioner summarized the events in the video as a video description.

718 Table 7 shows an example of a dialogue, and Table 8 shows the size of
719 the dataset split into training, validation, and test sets. The questions and
720 answers of the AVSD dataset mainly consist of 5 to 8 words, making them
721 longer and more descriptive than those of VQA and Visual Dialog. Figure
722 4 shows the distributions of word 4-grams and average length of sentences
723 in the questions and answers of the prototype data set of AVSD Hori et al.
724 (2019a), compared with those of VQA and Visual Dialog (VisDial).

725 The dialog contains questions about objects, actions, and audio informa-
726 tion in the videos. Although we tried to collect questions directly relevant
727 to the event displayed, some questions refer to abstract information in the
728 video, such as how the videos begin and the duration of the videos.

729 *4.3. Evaluation*

730 In this challenge, the quality of a system’s automatically generated sen-
731 tences is evaluated using objective measures. These determine how similar
732 the generated responses are to groundtruth responses from humans, as well
733 as how natural and informative the responses are. In addition to the ground
734 truth response that was given by the answerer during dialog collection, we
735 collected 5 additional human-generated responses for the test videos. To
736 collect these additional responses, we provided 5 humans with all of the in-
737 formation that the answerer had in the original dialog: we asked them to

¹³<http://allenai.org/plato/charades/>

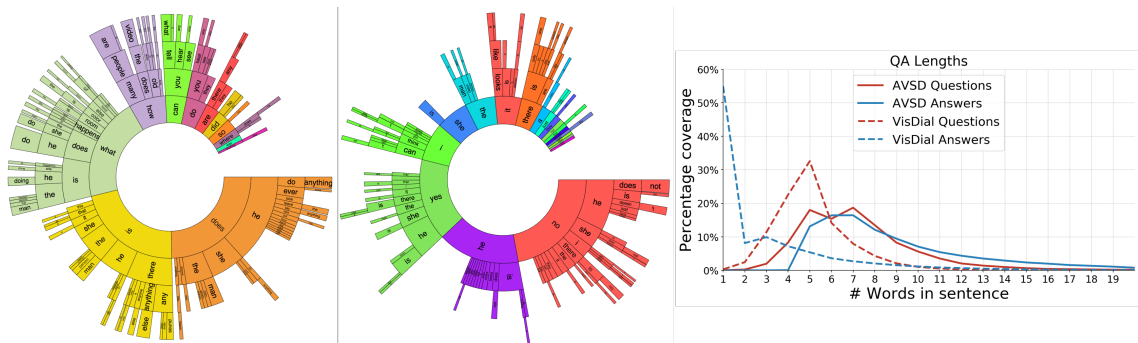


Figure 4: The distributions of word 4-grams in the questions (left) and answers (middle) of the prototype data set of the AVSD, and the average length (right) of the sentences of the VQA and the prototype data set of the AVSD. The actions were mainly asked by the questioners. There are some questions regarding audio information. Half of the answers are Yes/No. The questions and answers of AVSD are longer than those of VQA. More descriptive sentences were generated for AVSD.

738 answer the question after watching a video and reading the video script and
739 the dialog history between the questioner and answerer about the video. The
740 reason why the humans need to read the history of the dialog before answer-
741 ing is that there are some dependencies between each question and the the
742 previous question/answer pairs in the sequence (Alamri et al., 2019). A typi-
743 cal pattern is when questions contain prepositions such as "it" — the humans
744 cannot answer the questions if they don't know what the word "it" refers to.

745 We evaluated the automatically generated answers by comparing with the
746 6 ground truth sentences (one original answer and 5 subsequently collected
747 answers). We used the MSCOCO evaluation tool for objective evaluation of
748 system outputs¹⁴. The supported metrics include word-overlap-based metrics
749 such as BLEU, METEOR, ROUGE_L, and CIDEr.

750 We also collected human ratings for each system response using a 5-point
751 Likert Scale, where humans rated system responses given a dialog context as:
752 5 for very good, 4 for good, 3 for acceptable, 2 for poor, and 1 for very poor.
753 Since the dataset contains questions and answers, we asked humans to con-
754 sider correctness of the answers as well as the naturalness, informativeness,
755 and appropriateness of the response according to the given context.

¹⁴<https://github.com/tylin/coco-caption>

756 4.4. Baseline System

757 We provided a baseline end-to-end dialog system that can generate answers
 758 in response to user questions about events in a video sequence. The
 759 baseline system is an LSTM-based encoder decoder with Naïve multimodal
 760 fusion (Alamri et al., 2018b). The architecture, which is similar to the Hier-
 761 archical Recurrent Encoder in Das et al. (2016), is based on Natural language
 762 Generation (NLG) technologies from Track2 of DSTC6 (modeling end-to-end
 763 conversation for Twitter customer service) (Hori et al., 2018b). The question,
 764 visual features, and dialog history are fed into corresponding LSTM-based
 765 encoders to build up a context embedding, and then the outputs of the encoders
 766 are fed into an LSTM-based decoder to generate an answer. The
 767 dialog history consists of encodings of QA pairs plus (optionally) an encoding
 768 of the video script. This is a simplified version of Hori et al. (2018a), in
 769 which multimodal fusion is performed without attention between modalities
 770 such as audio and video features. Figure 5 shows the architecture of the mul-
 771 timodal attention-based fusion. The baseline system does not have modality
 772 attention weights β . The full set of test data was used in Hori et al. (2018a),
 773 while the AVSD challenge at DSTC7 used 2,000 responses selected from the
 774 full set.

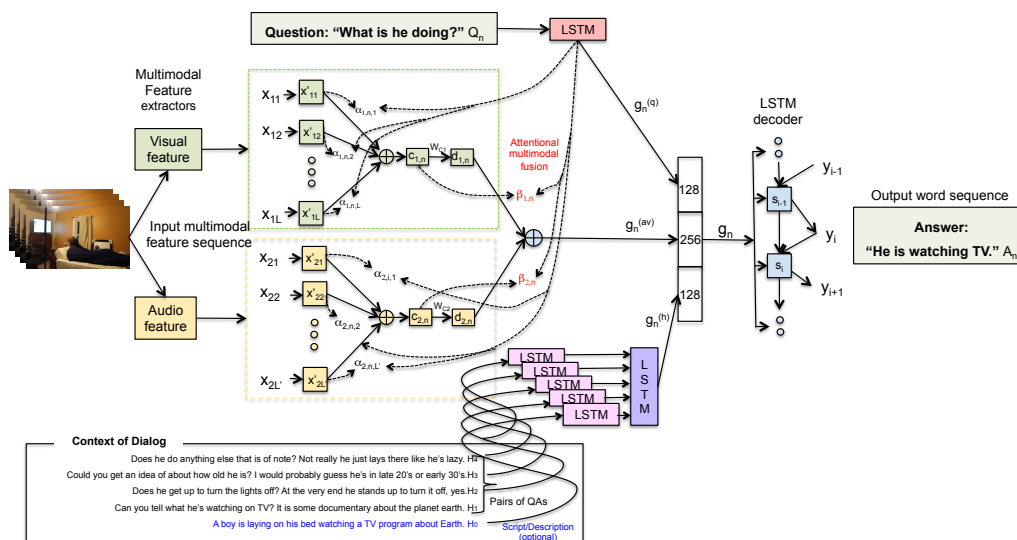


Figure 5: Attentional multimodal fusion-based video scene-aware dialog system Hori et al. (2018a)

775 *4.5. Data Processing*

776 *4.5.1. Video Processing*

777 We adopted the state-of-the-art I3D features Carreira and Zisserman
778 (2017), spatiotemporal features that were developed for action recognition.
779 The I3D model inflates the 2D filters and pooling kernels in the Inception V3
780 network along their temporal dimension, building 3D spatiotemporal ones.
781 We used the output from the "Mixed_5c" layer of the I3D network to be
782 used as video features in our framework. As a pre-processing step, we nor-
783 malized all the video features to have zero mean and unit norm; the mean
784 was computed over all the sequences in the training set for the respective
785 feature.

786 In the experiments in this paper, we treated I3D-rgb (I3D features com-
787 puted on a stack of 16 video frame images) and I3D-flow (I3D features com-
788 puted on a stack of 16 frames of optical flow fields) as two separate modalities
789 that are input to our multimodal attention model. To emphasize this, we
790 refer to I3D in the results tables as I3D (rgb-flow).

791 *4.5.2. Audio Processing*

792 In this track, we used features extracted using a new state-of-the-art
793 model, Audio Set VGGish (Hershey et al., 2017). Inspired by the VGG
794 image classification architecture (Configuration A without the last group of
795 convolutional/pooling layers), the Audio Set VGGish model operates on 0.96
796 sec log Mel spectrogram patches extracted from 16 kHz audio, and outputs
797 a 128-dimensional embedding vector. The model was trained to predict an
798 ontology of labels from only the audio tracks of millions of YouTube videos.
799 In this work, we overlap frames of input to the VGGish network by 50%,
800 meaning an Audio Set VGGish feature vector is output every 0.48 sec.

801 *4.6. Submitted Systems*

802 We received 32 sets of system outputs for the AVSD task, from 9 teams,
803 and eight system description papers were accepted (Sanabria et al., 2019;
804 Nguyen et al., 2019; Pasunuru and Bansal, 2019; Yeh et al., 2019; Zhuang
805 et al., 2019; Kumar et al., 2019; Lin et al., 2019; Le et al., 2019). Table 9 shows
806 the baseline and submitted systems with their brief specifications including
807 Encoder-decoder Model type, Multimodal fusion type, and Additional tech-
808 niques, models, and data sets. Most systems employed an LSTM, Bi-LSTM,
809 or GRU encoder/decoder. Some systems used hierarchical and attention

810 frameworks. Furthermore, several additional techniques were introduced to
 811 improve the response quality, such as MMI and Episodic Memory Module.

812 4.7. Results

813 The best system applied “Hierarchical Attention mechanisms to combine
 814 text and video,” which was proposed in Hori et al. (2018a). Table 10 shows
 815 the evaluation results for the baseline and all systems. Figures 6–8 show the
 816 human ratings for each system in several ways. The systems are shown in
 817 the same order on the x -axis for all three figures. Figure 6 shows the mean
 818 and the standard deviation of the human ratings for each system (across all
 819 responses and all raters for that system). Figure 7 shows the distributions
 820 of the mean human rating score for each sentence for each system. Figure 8
 821 shows the distribution of all human rating scores for each system across all
 822 sentences. In this Figure, the area for each score of the violin plot shows a
 823 count of the number of scores of each level on the Likert scale. The ratings of
 824 the reference system (labeled “Ref,” at the far left of each figure) are ratings
 825 for the ground truth sentences extracted from the original QA data of the
 826 AVSD dataset. The baseline system is labeled “Base.” The Reference system
 827 (“Ref”) had the best human ratings: it had the highest mean rating in Fig. 6,
 828 the highest median sentence rating in Fig. 7 and the most sentences rated as
 829 level 5 (“Very good”) in Fig. 8. The worst system (at the right) had a much
 830 lower mean rating and a long tail of poorly rated sentences.

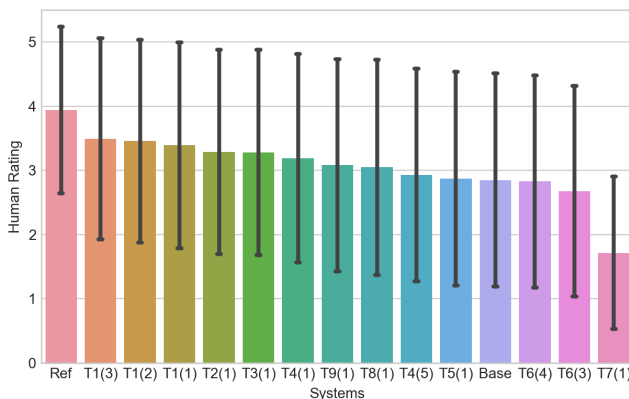


Figure 6: Mean and standard deviation of human rating score.

831 In Hori et al. (2018b), the reported human ratings of end-to-end con-
 832 versation models for Twitter customer service data were distributed fairly

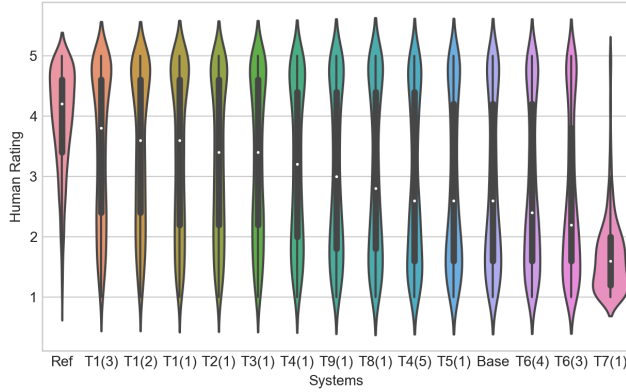


Figure 7: Distribution of human scores averaged sentence-by-sentence.

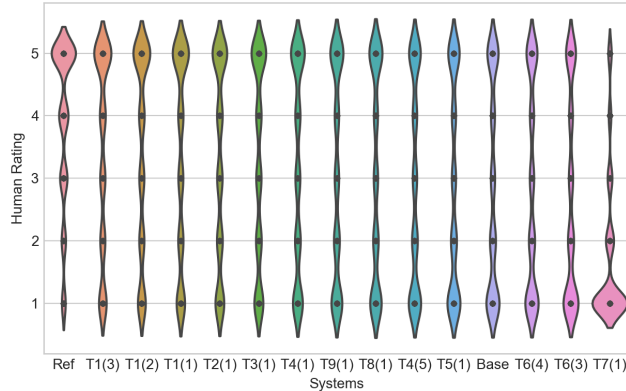


Figure 8: Distribution of human rating score for each level of scores.

833 smoothly in the range from 1 to 5. In contrast, the human ratings of re-
 834 sponses in this AVSD track were more bimodal, tending to be either very
 835 low or very high (more like a binary split into “good” and “bad” answers).
 836 This is because the quality of the answers depends on the answer correctness
 837 in response to the questions, and incorrect answers result in drastically lower
 838 human rating scores. The best system generated mostly correct answers, and
 839 the worst system generated mostly incorrect answers.

840 4.8. Summary and Discussion

841 We introduced a new challenge task and dataset for Audio Visual Scene-
 842 Aware Dialog (AVSD) in DSTC7. This is the first attempt to combine end-to-
 843 end conversation and end-to-end multimodal video description models into a
 844 single end-to-end differentiable network to build scene-aware dialog systems.

845 The best system applied hierarchical attention mechanisms to combine text
846 and visual information, improving by 22% over the human ratings of the
847 baseline system. The language models trained from QA (without video or
848 audio) are still strong approaches.

849 After the AVSD challenge at DSTC7, Alamri et al. (2019) reported the
850 performance of sentence selection (as opposed to sentence generation, which
851 was used in this AVSD challenge) using the AVSD dataset. In the paper,
852 **Question (Q)**, **V (Video)**, **Dialog History (DH)**, and **Audio (A)** were fused.
853 The addition of audio features generally improves model performance (**Q+V** to
854 **Q+V+A** being the exception). Interestingly, the model performance improves
855 even more when combined with dialog history and video features (**Q+DH+V+A**)
856 for some metrics, indicating that audio signals still provide complementary
857 knowledge to the video signals despite their close relationship.

858 Further, it is found that the best performance is achieved when including
859 text features extracted from the available summary (video script). Surpris-
860 ingly, systems that use such manual descriptions enable performance close
861 to the best system, even without using the audio-visual features. However,
862 such summaries are unavailable in the real world, posing challenges during
863 deployment. Recently, Hori et al. (2019b) proposed an approach to transfer
864 the power of the teacher model trained using summaries to a student model
865 that does not need the summary features.

866 5. Conclusion and Future Directions

867 In this paper, we have described the seventh dialog system technology
868 challenge (DSTC7) and the three selected tasks: sentence selection, sentence
869 generation, and audio visual scene-aware dialog. The sentence selection track
870 targeted the process of determining the best response given several possible
871 answers or detecting when none candidate was suitable over two different
872 datasets. The sentence generation track provided a testbed for knowledge-
873 grounded response generation, with the aim of creating more controllable
874 generators. The audio visual scene-aware dialog track raised a new prob-
875 lem in which dialog is generated about a given video, targeting multimodal
876 approaches and extending the capabilities of the dialog systems to combine
877 information from different sources.

878 All of the data described in this paper are provided as a large-scale bench-
879 mark of dialog systems from several viewpoints to support future dialog sys-
880 tem research. Although submitted systems improved in all cases the base-

881 line results, several major challenges for dialog systems still remain. For
882 example, transferring models trained on large-scale data-sets to a variety
883 of domains that do not have enough data is a known issue for dialog sys-
884 tems, as mentioned in DSTC3. Unfortunately, end-to-end systems do not
885 address completely this issue, which would require expanding to a larger
886 variety of domains and to consider applying transfer-learning approaches
887 (Ruder et al. (2019)). Other problems are related with the capabilities of the
888 dialog systems is to identify success and better managing of errors, handle
889 task complexity in a scalable way, and the integration of multiple sources of
890 information.

891 As following the raised problems in DSTC7, four tasks are proposed as
892 the eighth edition of the dialog system technology challenge (DSTC8). Sen-
893 tence selection task, track 1 in DSTC7, was extended not only a next ut-
894 terance selection task but also predicting a task success and a conversation
895 disentanglement. Audio visual scene aware dialog, track 3 in DSTC7, was
896 also continued in the next challenge to explore a fusion between vision and
897 dialog. Other two tasks, multi-domain task completion and scheme based di-
898 alog state tracking, were proposed as new challenges in DSTC8. Both tracks
899 aim to build accurate task-oriented dialog systems on different approaches.
900 Multi-domain task completion track focuses on dialog complexity and scaling
901 to new domains as we previously focused on DSTC3. Scheme guided dialog
902 state tracking focuses on dialog state tracking itself, even if the state space
903 is new for the trained state tracker.

904 We expect to continue the challenge in the future, providing new testbeds
905 that work towards the remaining open problems of dialog system research,
906 while being complementary to other challenges like Alexa Prize (Khatri et al.
907 (2018)), ConvAI (Dinan et al. (2019)), or Dialog Breakdown Detection Chal-
908 lenge (Higashinaka et al. (2019)).

909 **6. Bibliography**

910 Alamri, H., Cartillier, V., Das, A., Wang, J., Cherian, A., Essa, I., Batra,
911 D., Marks, T.K., Hori, C., Anderson, P., Lee, S., Parikh, D., 2019. Audio
912 visual scene-aware dialog, in: The IEEE Conference on Computer Vision
913 and Pattern Recognition (CVPR).

914 Alamri, H., Cartillier, V., Lopes, R.G., Das, A., Wang, J., Essa, I., Batra,
915 D., Parikh, D., Cherian, A., Marks, T.K., et al., 2018a. Audio visual scene-

- 916 aware dialog (avsd) challenge at dstc7. arXiv preprint arXiv:1806.00525
917 .
- 918 Alamri, H., Hori, C., Marks, T.K., Batra, D., Parikh, D., 2018b. Audio
919 visual scene-aware dialog (avsd) track for natural language generation in
920 dstc7, in: DSTC7 at AAIL2019 Workshop.
- 921 Antol, S., Agrawal, A., Lu, J., Mitchell, M., Batra, D., Zitnick, C.L., Parikh,
922 D., 2015. VQA: Visual Question Answering, in: International Conference
923 on Computer Vision (ICCV).
- 924 Bahdanau, D., Cho, K., Bengio, Y., 2015. Neural machine translation by
925 jointly learning to align and translate, in: Proc. of the International Con-
926 ference on Learning Representations (ICLR).
- 927 Banchs, R.E., D’Haro, L.F., Li, H., 2015. Adequacy–fluency metrics: Eval-
928 uating mt in the continuous space model framework. *IEEE/ACM Trans-*
929 *actions on Audio, Speech, and Language Processing* 23, 472–482.
- 930 Carreira, J., Zisserman, A., 2017. Quo vadis, action recognition? a new
931 model and the kinetics dataset, in: CVPR.
- 932 Chen, Q., Zhu, X., Ling, Z.H., Wei, S., Jiang, H., Inkpen, D.,
933 2017. Enhanced LSTM for natural language inference, in: Proceed-
934 ings of the 55th Annual Meeting of the Association for Computa-
935 tional Linguistics (Volume 1: Long Papers), pp. 1657–1668. URL:
936 <http://aclweb.org/anthology/P17-1152>, doi:10.18653/v1/P17-1152.
- 937 Chen, Q.Q., Wang, W., 2019. Sequential attention-based network for noetic
938 end-to-end response selection, in: 7th Edition of the Dialog System Tech-
939 nology Challenges at AAIL 2019.
- 940 Das, A., Kottur, S., Gupta, K., Singh, A., Yadav, D., Moura, J.M.F.,
941 Parikh, D., Batra, D., 2016. Visual dialog. CoRR abs/1611.08669. URL:
942 <http://arxiv.org/abs/1611.08669>, arXiv:1611.08669.
- 943 Das, A., Kottur, S., Moura, J.M., Lee, S., Batra, D., 2017. Learning coopera-
944 tive visual dialog agents with deep reinforcement learning, in: International
945 Conference on Computer Vision (ICCV).

- 946 Devlin, J., Chang, M.W., Lee, K., Toutanova, K., 2018. Bert: Pre-training of
947 deep bidirectional transformers for language understanding. arXiv preprint
948 arXiv:1810.04805 .
- 949 D’Haro, L.F., Banchs, R.E., Hori, C., Li, H., 2019. Automatic evaluation
950 of end-to-end dialog systems with adequacy-fluency metrics. *Computer
951 Speech & Language* 55, 200–215.
- 952 Dinan, E., Logacheva, V., Malykh, V., Miller, A., Shuster, K., Urbanek, J.,
953 Kiela, D., Szlam, A., Serban, I., Lowe, R., et al., 2019. The second conver-
954 sational intelligence challenge (convai2). arXiv preprint arXiv:1902.00098
955 .
- 956 Doddington, G., 2002. Automatic evaluation of machine translation qual-
957 ity using n-gram co-occurrence statistics, in: *Proceedings of the Second
958 International Conference on Human Language Technology Research*, pp.
959 138–145.
- 960 Ganhotra, J., Patel, S.S., Fadnis, K.P., 2019. Knowledge-incorporating ESIM
961 models for response selection in retrieval-based dialog systems, in: *7th
962 Edition of the Dialog System Technology Challenges at AAAI 2019*.
- 963 Gao, X., Lee, S., Zhang, Y., Brockett, C., Galley, M., Gao, J., Dolan, B.,
964 2019. Jointly optimizing diversity and relevance in neural response gener-
965 ation. *Proceedings of the 2019 Conference of the North American Chapter
966 of the Association for Computational Linguistics* .
- 967 Ghazvininejad, M., Brockett, C., Chang, M., Dolan, B., Gao, J., Yih, W.,
968 Galley, M., 2018. A knowledge-grounded neural conversation model. *AAAI
969* .
- 970 Gu, J., Lu, Z., Li, H., Li, V.O., 2016. Incorporating copying mechanism in
971 sequence-to-sequence learning, in: *Proceedings of the 54th Annual Meeting
972 of the Association for Computational Linguistics (Volume 1: Long Papers)*.
- 973 He, S., Liu, C., Liu, K., Zhao, J., 2017. Generating natural answers by
974 incorporating copying and retrieving mechanisms in sequence-to-sequence
975 learning, in: *ACL*, pp. 199–208.
- 976 Henderson, M., Thomson, B., Williams, J.D., 2014a. The second dialog
977 state tracking challenge, in: *Proceedings of the 15th Annual Meeting of*

978 the Special Interest Group on Discourse and Dialogue (SIGDIAL), pp.
979 263–272.

980 Henderson, M., Thomson, B., Williams, J.D., 2014b. The third dialog state
981 tracking challenge, in: Spoken Language Technology Workshop (SLT),
982 2014 IEEE, IEEE. pp. 324–329.

983 Hershey, S., Chaudhuri, S., Ellis, D.P.W., Gemmeke, J.F., Jansen, A., Moore,
984 R.C., Plakal, M., Platt, D., Saurous, R.A., Seybold, B., Slaney, M., Weiss,
985 R.J., Wilson, K., 2017. CNN architectures for large-scale audio classifica-
986 tion, in: ICASSP.

987 Higashinaka, R., D’Haro, L.F., Shawar, B.A., Banchs, R., Funakoshi,
988 K., Inaba, M., Tsunomori, Y., Takahashi, T., Sedoc, J., 2019.
989 Overview of the dialogue breakdown detection challenge 4, in: 10th
990 International Workshop on Spoken Dialog Systems (IWSDS). URL:
991 <http://workshop.colips.org/wchat/@iwsds2019/documents/dbdc4-overview-higashinaka>

992 Hori, C., Alamri, H., Wang, J., Wichern, G., Hori, T., Cherian, A., Marks,
993 T.K., Cartillier, V., Lopes, R.G., Das, A., et al., 2019a. End-to-end au-
994 dio visual scene-aware dialog using multimodal attention-based video fea-
995 tures, in: ICASSP 2019-2019 IEEE International Conference on Acoustics,
996 Speech and Signal Processing (ICASSP), IEEE. pp. 2352–2356.

997 Hori, C., Alamri, H., Wang, J., Winchern, G., Hori, T., Cherian, A., Marks,
998 T.K., Cartillier, V., Lopes, R.G., Das, A., et al., 2018a. End-to-end audio
999 visual scene-aware dialog using multimodal attention-based video features.
1000 arXiv preprint arXiv:1806.08409 .

1001 Hori, C., Hori, T., 2017a. End-to-end conversation modeling track
1002 in DSTC6, in: Dialog System Technology Challenges 6. URL:
1003 http://workshop.colips.org/dstc6/papers/track2_overview_hori.pdf.

1004 Hori, C., Hori, T., 2017b. End-to-end conversation modeling track in DSTC6.
1005 arXiv:1706.07440 .

1006 Hori, C., Hori, T., Cherian, A., Marks, T.K., 2019b. Joint student-teacher
1007 learning for audio-visual scene-aware dialog, in: Interspeech2019, ISCA.

- 1008 Hori, C., Hori, T., Lee, T.Y., Zhang, Z., Harsham, B., Hershey, J.R., Marks,
1009 T.K., Sumi, K., 2017. Attention-based multimodal fusion for video de-
1010 scription, in: ICCV.
- 1011 Hori, C., Perez, J., Higashinaka, R., Hori, T., Boureau, Y.L., Inaba, M.,
1012 Tsunomori, Y., Takahashi, T., Yoshino, K., Kim, S., 2019c. Overview of
1013 the sixth dialog system technology challenge: Dstc6. *Computer Speech &
1014 Language* 55, 1–25.
- 1015 Hori, C., Perez, J., Higashinaka, R., Hori, T., Boureau, Y.L., Inaba, M.,
1016 Tsunomori, Y., Takahashi, T., Yoshino, K., Kim, S., 2018b. Overview of
1017 the sixth dialog system technology challenge: DSTC6. *Computer Speech
1018 and Language Special issue on DSTC6*.
- 1019 Jiang, Y., Kummerfeld, J.K., Lasecki, W.S., 2017. Understand-
1020 ing task design trade-offs in crowdsourced paraphrase collection,
1021 in: *Proceedings of the 55th Annual Meeting of the Association
1022 for Computational Linguistics (Volume 2: Short Papers)*. URL:
1023 <http://aclweb.org/anthology/P17-2017>.
- 1024 Khatri, C., Hedayatnia, B., Venkatesh, A., Nunn, J., Pan, Y., Liu, Q., Song,
1025 H., Gottardi, A., Kwatra, S., Pancholi, S., et al., 2018. Advancing the
1026 state of the art in open domain dialog systems through the alexa prize.
1027 arXiv preprint arXiv:1812.10757 .
- 1028 Kim, S., D’Haro, L.F., Banchs, R.E., Williams, J.D., Henderson, M., 2017.
1029 The fourth dialog state tracking challenge, in: *Dialogues with Social
1030 Robots*. Springer, pp. 435–449.
- 1031 Kim, S., D’Haro, L.F., Banchs, R.E., Williams, J.D., Henderson, M.,
1032 Yoshino, K., 2016. The fifth dialog state tracking challenge, in: *Spoken
1033 Language Technology Workshop (SLT), 2016 IEEEover, IEEE*. pp.
1034 511–517.
- 1035 Kingma, D.P., Welling, M., 2013. Auto-encoding variational bayes. arXiv
1036 preprint arXiv:1312.6114 .
- 1037 Kumar, S.H., Okur, E., Sahay, S., Leanos, J.J.A., Huang, J., Nachman, L.,
1038 2019. Context, attention and audio feature explorations for audio visual
1039 scene-aware dialogue, in: *DSTC7 at AAI2019 workshop*.

- 1040 Kummerfeld, J.K., 2019. Slate: A super-lightweight annotation tool for ex-
1041 perts, in: Proceedings of ACL 2019, System Demonstrations.
- 1042 Kummerfeld, J.K., Gouravajhala, S.R., Peper, J., Athreya, V., Gu-
1043 nasekara, C., Ganhotra, J., Patel, S.S., Polymenakos, L., Lasecki, W.S.,
1044 2018. Analyzing assumptions in conversation disentanglement research
1045 through the lens of a new dataset and model. ArXiv e-prints URL:
1046 <https://arxiv.org/pdf/1810.11118.pdf>, arXiv:1810.11118.
- 1047 Lavie, A., Agarwal, A., 2007. METEOR: An automatic metric for mt eval-
1048 uation with high levels of correlation with human judgments, in: Proc.
1049 of the Second Workshop on Statistical Machine Translation, Association
1050 for Computational Linguistics, Stroudsburg, PA, USA. pp. 228–231. URL:
1051 <http://dl.acm.org/citation.cfm?id=1626355.1626389>.
- 1052 Le, H., Hoi, S., Sahoo, D., Chen, N., 2019. End-to-end multimodal dia-
1053 log systems with hierarchical multimodal attention on video features, in:
1054 DSTC7 at AAIL2019 workshop.
- 1055 Li, J., Galley, M., Brockett, C., Gao, J., Dolan, B., 2016a. A diversity-
1056 promoting objective function for neural conversation models, in: NAACL-
1057 HLT.
- 1058 Li, J., Galley, M., Brockett, C., Gao, J., Dolan, B., 2016b. A diversity-
1059 promoting objective function for neural conversation models, in: Proceed-
1060 ings of the 2016 Conference of the North American Chapter of the Associ-
1061 ation for Computational Linguistics: Human Language Technologies, pp.
1062 110–119.
- 1063 Li, J., Galley, M., Brockett, C., Gao, J., Dolan, B., 2016c. A diversity-
1064 promoting objective function for neural conversation models. NAACL-HLT
1065 .
- 1066 Lin, K.Y., Hsu, C.C., Chen, Y.N., Ku, L.W., 2019. Entropy-enhanced mul-
1067 timodal attention model for scene-aware dialogue generation, in: DSTC7
1068 at AAIL2019 workshop.
- 1069 Lowe, R., Pow, N., Serban, I., Pineau, J., 2015. The ubuntu dialogue
1070 corpus: A large dataset for research in unstructured multi-turn dia-
1071 logue systems, in: Proceedings of the 16th Annual Meeting of the Spe-

- 1072 cial Interest Group on Discourse and Dialogue, Association for Com-
1073 putational Linguistics, Prague, Czech Republic. pp. 285–294. URL:
1074 <http://aclweb.org/anthology/W15-4640>.
- 1075 Nguyen, D., Sharma, S., Schulz, H., Asri, L.E., 2019. From film to video:
1076 Multi-turn question answering with multi-modal context, in: DSTC7 at
1077 AAAI2019 workshop.
- 1078 Papineni, K., Roukos, S., Ward, T., Zhu, W.J., 2002. BLEU: a method for
1079 automatic evaluation of machine translation. ACL .
- 1080 Pasunuru, R.R., Bansal, M., 2019. Dstc7-avsd: Scene-aware video-dialogue
1081 systems with dual attention, in: DSTC7 at AAAI2019 workshop.
- 1082 Perez, J., Boureau, Y.L., Bordes, A., 2017. Dialog system technol-
1083 ogy challenge 6 overview of track 1 - end-to-end goal-oriented di-
1084 alog learning, in: Dialog System Technology Challenges 6. URL:
1085 http://workshop.colips.org/dstc6/papers/track1_overview_perez.pdf.
- 1086 Peters, M., Neumann, M., Iyyer, M., Gardner, M., Clark, C., Lee,
1087 K., Zettlemoyer, L., 2018. Deep contextualized word representations,
1088 in: Proceedings of the 2018 Conference of the North American Chap-
1089 ter of the Association for Computational Linguistics: Human Lan-
1090 guage Technologies, Volume 1 (Long Papers), pp. 2227–2237. URL:
1091 <http://aclweb.org/anthology/N18-1202>, doi:10.18653/v1/N18-1202.
- 1092 Qin, L., Galley, M., Brockett, C., Liu, X., Gao, X., Dolan, B., Choi, Y., Gao,
1093 J., 2019. Conversing by reading: Contentful neural conversation with on-
1094 demand machine reading, in: Proc. of ACL.
- 1095 Ritter, A., Cherry, C., Dolan, W.B., 2011. Data-driven response generation
1096 in social media. EMNLP .
- 1097 Ruder, S., Peters, M.E., Swayamdipta, S., Wolf, T., 2019. Transfer learning in
1098 natural language processing, in: Proceedings of the 2019 Conference of the
1099 North American Chapter of the Association for Computational Linguistics:
1100 Tutorials, pp. 15–18.
- 1101 Sanabria, R., Palaskar, S., Metze, F., 2019. Cmu sinbad submission for the
1102 dstc7 avsd challenge, in: DSTC7 at AAAI2019 workshop.

- 1103 See, A., Liu, P.J., Manning, C.D., 2017. Get to the point: Summariza-
1104 tion with pointer-generator networks, in: Proceedings of the 55th Annual
1105 Meeting of the Association for Computational Linguistics (Volume 1: Long
1106 Papers), pp. 1073–1083. doi:10.18653/v1/P17-1099.
- 1107 Serban, I.V., Lowe, R., Henderson, P., Charlin, L., Pineau, J., 2018.
1108 A survey of available corpora for building data-driven dialogue sys-
1109 tems: The journal version. *Dialogue & Discourse* 9, 1–49. URL:
1110 <http://dad.uni-bielefeld.de/index.php/dad/article/view/3690>,
1111 doi:10.5087/dad.2018.101.
- 1112 Serban, I.V., Sordoni, A., Bengio, Y., Courville, A., Pineau, J., 2016. Build-
1113 ing end-to-end dialogue systems using generative hierarchical neural net-
1114 work models, in: Proc. of AAAI.
- 1115 Shang, L., Lu, Z., Li, H., 2015. Neural responding machine for short-text
1116 conversation. *ACL-IJCNLP* .
- 1117 Sharma, S., El Asri, L., Schulz, H., Zumer, J., 2017. Rele-
1118 vance of unsupervised metrics in task-oriented dialogue for evalu-
1119 ating natural language generation. *CoRR* abs/1706.09799. URL:
1120 <http://arxiv.org/abs/1706.09799>.
- 1121 Sigurdsson, G.A., Varol, G., Wang, X., Laptev, I., Farhadi, A., Gupta,
1122 A., 2016. Hollywood in homes: Crowdsourcing data collection for ac-
1123 tivity understanding. *ArXiv* URL: <http://arxiv.org/abs/1604.01753>,
1124 arXiv:1604.01753.
- 1125 Sordoni, A., Galley, M., Auli, M., Brockett, C., Ji, Y., Mitchell, M., Nie, J.Y.,
1126 Gao, J., Dolan, B., 2015. A neural network approach to context-sensitive
1127 generation of conversational responses. *NAACL-HLT* .
- 1128 Sukhbaatar, S., szlam, a., Weston, J., Fergus, R., 2015. End-to-end mem-
1129 ory networks, in: *Advances in Neural Information Processing Systems 28*.
1130 Curran Associates, Inc., pp. 2440–2448.
- 1131 Vedantam, R., Zitnick, C.L., Parikh, D., 2015a. CIDEr: Consensus-based
1132 image description evaluation., in: *CVPR*, pp. 4566–4575.
- 1133 Vedantam, R., Zitnick, C.L., Parikh, D., 2015b. CIDEr: Consensus-based
1134 image description evaluation, in: *IEEE Conference on Computer Vision*

- 1135 and Pattern Recognition, CVPR 2015, Boston, MA, USA, June 7-12, 2015,
1136 pp. 4566–4575.
- 1137 Vinyals, O., Le, Q., 2015. A neural conversational model. ICML .
- 1138 Weston, J., Chopra, S., Bordes, A., 2015. Memory networks. ICLR .
- 1139 Williams, J., Raux, A., Ramachandran, D., Black, A., 2013. The dialog state
1140 tracking challenge, in: Proceedings of the SIGDIAL 2013 Conference, pp.
1141 404–413.
- 1142 Yeh, Y.T., Lin, T.C., Cheng, H.H., Deng, Y.H., Su, S.Y., Chen, Y.N., 2019.
1143 Reactive multi-stage feature fusion for multimodal dialogue modeling, in:
1144 DSTC7 at AAI2019 workshop.
- 1145 Zhuang, B., Wang, W., Shinozaki, T., 2019. Investigation of attention-based
1146 multimodal fusion and maximum mutual information objective for dstc7
1147 track3, in: DSTC7 at AAI2019 workshop.

Team	Model Type	External Data Use	Used Raw Advising	Val in Train	Model Details
1	CNN	-	No	Yes	Combination of CNN for utterance representation and GRU for modeling the dialogue.
2	LSTM	-	Yes	No	ESIM with an aggregation scheme to capture dialog-specific aspects of the data + ELMo.
3	LSTM	Embeddings	Yes	No	ESIM + a filtering stage for subtask 2.
4	LSTM	-	No	No	ESIM with (1) enhanced word embeddings to address OOV issues, (2) an attentive hierarchical recurrent encoder, and (3) an additional layer before the softmax.
6	Ensemble	-	No	No	An ensemble of CNNs.
7	LSTM	-	No	Yes	LSTM representation of utterances followed by a convolutional layer.
8	Other	-	Yes	No	A multi-level retrieval-based approach that aggregates similarity measures between the context and the candidate response on the sequence and word levels.
10	LSTM	TF-IDF Extraction	No	No	ESIM with matching against similar dialogues in training, and an extra filtering step for subtask 2.
12	RNN	TF-IDF Extraction	No	No	BoW over ELMo with context as an RNN.
13	Ensemble	Embeddings	No	No	Ensemble approach, combining a Dynamic-Pooling LSTM, a Recurrent Transformer and a Hierarchical LSTM.
14	Ensemble	-	No	No	An ensemble using voting, combining the baseline LSTM, a GRU variant, Doc2Vec, TF-IDF, and LSI.
15	Memory	Memory	No	No	Memory network with an LSTM cell.
16	LSTM	-	No	No	ESIM with utterance-level attention, plus additional features.
17	Memory	Memory & Embeddings	Yes	No	Self-attentive memory network, with external advising data in memory and external ubuntu data for embedding training.
18	GRU	-	No	No	Stacked Bi-GRU network with attention, aggregating attention across the temporal dimension followed by a CNN and softmax.
19	LSTM	-	No	Yes	Bidirectional LSTM memory network.
20	CNN	-	No	Yes	CNN with attention and a pointer network, plus a novel top-k attention mechanism.

Table 2: Summary of approaches used by participants for track-1. All teams applied neural approaches, with ESIM being a popular basis for system development. External data refers to the man pages for Ubuntu, and course information for Advising. Raw advising refers to the variant of the training data in which the complete dialogues and paraphrase sets are provided. Teams 5, 9 and 11 did not provide descriptions of their approaches. For further details, see the system description papers presented at the DSTC workshop.

Team	Ubuntu, Subtask				Advising, Subtask			
	1	2	4	5	1	3	4	5
3	0.819	0.145	0.842	0.822	0.485	0.592	0.537	0.485
4	0.772	-	-	-	0.451	-	-	-
17	0.705	-	-	0.722	0.434	-	-	0.461
13	0.729	-	0.736	0.635	0.458	0.461	0.474	0.390
2	0.672	0.033	0.713	0.672	0.430	0.540	0.479	0.430
10	0.651	0.307	0.696	0.693	0.361	0.434	0.262	0.361
18	0.690	0.000	0.721	0.710	0.287	0.380	0.398	0.326
8	0.641	-	0.527	-	0.310	0.433	0.233	-
16	0.629	0.000	0.683	-	0.280	-	0.370	-
15	0.473	-	-	0.478	0.300	-	-	0.236
7	0.525	-	0.411	-	-	-	-	-
11	-	-	-	-	0.075	0.232	-	-
12	0.077	-	0.000	0.077	0.075	0.232	0.000	0.075
1	0.580	-	-	-	0.239	-	-	-
6	-	-	-	-	0.245	-	-	-
9	0.482	-	-	-	-	-	-	-
14	0.008	-	0.072	-	-	-	-	-
19	0.265	-	-	-	0.180	-	-	-
5	0.076	-	-	-	-	-	-	-
20	0.002	-	-	-	0.004	-	-	-

Table 3: Track-1 results, ordered by the average rank of each team across the sub-tasks they participated in. The top result in each column is in bold. For these results the metric is the average of MRR and Recall@10.

<i>Web page info</i>	[...] she holds the guinness world record for surviving the highest fall without a parachute : 10,160 metres (33,330 ft) . [...] four years later , peter hornung-andersen and pavel theiner , two prague-based journalists , claimed that flight 367 had been mistaken for an enemy aircraft and shot down by the czechoslovak air force at an altitude of 800 metres (2,600 ft) [...]
<i>Turn 1</i>	today i learned a woman fell 30,000 feet from an airplane and survived [URL] .
<i>Turn 2</i>	the page states that a 2009 report found the plane only fell several hundred meters .
<i>Turn 3</i>	well if she only fell a few hundred meters and survived then i 'm not impressed at all .
<i>Turn 4</i>	still pretty incredible , but quite a bit different that 10,000 meters .

Table 4: Sample of the DSTC7 Sentence Generation data, which combines Reddit data (Turns 1-4) along with documents (extracted from Common Crawl) discussed in the conversations. The web page info was truncated for this figure to fit in a relatively small space. The **emphasis** was added by us. The [URL] links to the web page above.

System	NIST		BLEU(%)		METEOR	Diversity		Avg. len
	N-2	N-4	B-2	B-4		D-1	D-2	
<i>Baselines:</i>								
Constant	0.18	0.18	12.8	2.9	7.5	0.1	0.1	8.0
Random_Human	1.63	1.64	6.7	0.9	5.9	16.0	64.7	19.2
Seq2Seq	0.91	0.92	14.8	1.8	7.0	1.4	4.8	10.6
TeamA	0.75	0.75	11.8	1.5	5.6	9.6	27.6	10.5
TeamA-c1	0.83	0.83	11.5	1.4	5.7	12.2	30.2	10.9
TeamA-c2	1.12	1.12	9.5	0.8	5.5	9.7	31.9	12.0
TeamB	2.51	2.5	14.4	1.8	8.1	10.9	32.5	15.1
TeamB-c1	1.76	1.77	13.7	1.9	7.6	9.4	26.7	12.8
TeamC	1.51	1.51	10.9	1.3	6.4	5.3	17.1	12.7
TeamC-c1	2.11	2.12	9.9	1.3	6.8	3.8	12.4	16.4
TeamC-c2	1.19	1.20	11.6	1.7	6.2	5.5	16.9	11.7
TeamC-c3	1.73	1.74	8.8	1.2	5.9	3.9	12.2	14.9
TeamC-c4	1.53	1.54	11.5	1.8	6.5	5.6	18.0	12.7
TeamD	2.04	2.05	11.3	1.4	6.7	9.4	33.4	14.4
TeamD-c1	0.02	0.02	6.7	0.3	3.9	2.6	16.1	6.2
TeamD-c2	0.73	0.73	9.3	0.6	5.7	4.9	31.3	10.4
TeamD-c3	0.77	0.77	9.2	0.7	5.6	4.9	30.9	10.5
TeamD-c4	0.55	0.56	8.8	0.8	5.2	6.9	35.2	9.8
TeamD-c5	1.80	1.80	10.7	0.9	6.5	5.8	29.2	13.5
TeamD-c6	1.74	1.75	12.5	1.1	6.7	5.1	20.7	13.1
TeamE	1.51	1.51	10.9	1.3	6.4	5.3	17.1	12.7
TeamE-c1	2.11	2.12	9.9	1.3	6.8	3.8	12.4	16.4
TeamE-c2	1.81	1.82	11.0	1.6	6.5	5.0	15.6	14.0
TeamE-c3	1.92	1.93	10.9	1.5	6.7	4.6	15.2	14.3
TeamF	0.01	0.01	10.2	1.0	4.6	6.4	17.6	5.4
TeamF-c1	0.01	0.01	9.0	1.3	4.1	2.4	7.2	5.1
TeamF-c2	0.04	0.04	11.2	1.4	5.0	8.4	22.4	6.3
TeamG	2.31	2.32	10.6	1.2	7.2	3.4	26.5	16.6
TeamG-c1	2.03	2.04	8.2	1.1	7.5	10.8	44.9	22.3
Human	2.62	2.65	12.4	3.1	8.3	16.7	67.0	18.8

Table 5: Automatic evaluation results for track-2. Participants submitted primary and contrastive systems, the latter being identified with a -cX suffix in their names. The primary systems (TeamA, TeamB, ...) were the ones selected by the participants for human evaluation (Table 6).

System	Relevance		Interest		Overall	
	Mean	95% CI	Mean	95% CI	Mean	95 % CI
<i>Baselines:</i>						
Constant	2.60	(2.560, 2.644)	2.32	(2.281, 2.364)	2.46	(2.424, 2.500)
Random	2.32	(2.269, 2.371)	2.35	(2.303, 2.401)	2.34	(2.288, 2.384)
Seq2Seq	2.91	(2.858, 2.963)	2.68	(2.632, 2.730)	2.80	(2.748, 2.844)
TeamA	2.32	(2.267, 2.368)	2.30	(2.252, 2.351)	2.31	(2.262, 2.358)
TeamB	2.99	(2.938, 3.042)	2.87	(2.822, 2.922)	2.93	(2.882, 2.979)
TeamC	3.05	(3.009, 3.093)	2.77	(2.735, 2.812)	2.91	(2.875, 2.950)
TeamD	2.69	(2.635, 2.743)	2.58	(2.527, 2.632)	2.63	(2.583, 2.685)
TeamF	2.52	(2.461, 2.572)	2.40	(2.352, 2.457)	2.46	(2.409, 2.512)
TeamG	2.82	(2.771, 2.870)	2.57	(2.525, 2.619)	2.70	(2.650, 2.742)
Human	3.61	(3.554, 3.658)	3.49	(3.434, 3.539)	3.55	(3.497, 3.596)

Table 6: Human evaluation results for track-2. The systems evaluated here are the same as the primary systems in Table 5. Note that we do not report the results of TeamE as their primary system was identical to TeamC’s (due to miss-communication at submission time). The best system according to human evaluation (TeamB) also obtained the best NIST-4 and METEOR scores.

	Questioner	Answerer
QA1	What kind of room does this appear to be?	He appears to be in the bedroom.
QA2	How does the video begin?	By him entering the room.
QA3	Does he have anything in his hands?	He pick up a towel and folds it.
QA4	What does he do with it ?	He just folds them and leaves them on the chair.
QA5	What does he do next?	Nothing much except this activity.
QA6	Does he speak in the video?	No he did not speak at all.
QA7	Is there anyone else in room at all?	No he appears alone there.
QA8	Can you see or hear any pets in the video?	No pets to see in this clip.
QA9	Is there any noise in the video of importance?	Not any noise important there.
QA10	Are there any other actions in the video?	Nothing else important to know.

Table 7: An example dialog from the AVSD dataset.

	training	validation	test
# of dialogs	7,659	1,787	1,710
# of turns	153,180	35,740	13,490
# of words	1,450,754	339,006	110,252

Table 8: The dialog data for the DSTC7 AVSD track. The test videos for this challenge were selected from the official test data of the Charades challenge.

Team	Encoder-decoder type	Multimodal fusion type	Additional techniques/data
baseline	LSTM	Naïve fusion	
team_1	Bidirectional Gated Recurrent Units (GRU) based encode, Conditional Gated Recurrent Units (CGRU) based decoder	Hierarchical attention	ResNeXt, Transfer learning using How2 dataset
team_2	FiLM Attention Hierarchical Recurrent Encoder Decoder (FA-HRED), LSTM	Naïve fusion	FiLM
team_3	Dual attention LSTM encoder,	Cross-attention fusion	Similarity matrix
team_4	LSTM/GRU encoder, Top-down Attention LSTM/GRU decoder	Muti-stage fusion, 1x1 Convolution fusion, Multi-head Attention	
team_5	Bi-LSTM and LSTM encoder, LSTM decoder	Attentional multimodal fusion	MMI objective
team_6	LSTM encoder-decoder	Attentional multimodal fusion	Topic-base Conceptual model, ConvNet, ActMet
team_7	–	–	–
team_8	Bi-LSTM/LSTM encoder, Attention-based GRU encoder, LSTM decoder	Entropy-enhanced Dynamic Memory Network (DMN)	Episodic Memory Module
team_9	GRU encoder-decoder	Question-to-Caption/Multimodal attention	

⁺Team 7 did not submit a system description paper to the DSTC7 workshop.

Table 9: Submitted systems to the AVSD Track.

Team	Entry	text only	video	caption and/or summary	extra	prototype	Bleu_4	METEOR	ROUGE_L	CIDEr	Human rating
Team 1	(1)	✓		✓	✓		0.376	0.264	0.554	1.076	3.394
	(2)		✓	✓	✓		0.387	0.266	0.564	1.087	3.459
	(3)		✓	✓			0.394	0.267	0.563	1.094	3.491
	(4)	✓		✓			0.364	0.254	0.543	1.006	-
Team 2	(1)		✓	✓			0.360	0.249	0.544	0.997	3.288
	(2)	✓	✓	✓			0.323	0.231	0.510	0.843	
	(3)	✓		✓			0.343	0.243	0.536	0.920	
	(4)	✓		✓			0.340	0.228	0.518	0.851	
	(5)			✓		✓	0.349	0.242	0.536	0.947	
	(6)			✓	✓	✓	0.316	0.224	0.505	0.795	
	(7)			✓	✓	✓	0.319	0.228	0.513	0.836	
	(8)	✓		✓		✓	0.323	0.220	0.501	0.799	
Team 3	(1)		✓	✓			0.337	0.242	0.532	0.957	3.279
Team 4	(1)		✓	✓		✓	0.342	0.223	0.504	0.837	3.188
	(2)		✓	✓			0.345	0.224	0.505	0.877	
	(3)		✓	✓		✓	0.342	0.223	0.504	0.836	
	(4)	✓		✓			0.304	0.207	0.477	0.731	
	(5)	✓		✓			0.304	0.206	0.475	0.729	2.928
Team 5	(1)		✓			✓	0.293	0.221	0.486	0.761	2.869
	(2)		✓			✓	0.302	0.222	0.488	0.770	
	(3)		✓			✓	0.302	0.222	0.487	0.769	
	(4)		✓			✓	0.296	0.219	0.484	0.745	
	(5)		✓			✓	0.283	0.217	0.480	0.731	
Team 6	(1)	✓	✓	✓		✓	0.307	0.213	0.469	0.701	
	(2)	✓	✓	✓		✓	0.307	0.215	0.479	0.733	
	(3)	✓	✓	✓		✓	0.278	0.198	0.442	0.614	2.675
	(4)	✓	✓	✓		✓	0.310	0.217	0.483	0.718	2.827
Team 7	(1)	✓		✓			0.056	0.096	0.236	0.085	1.715
Team 8	(1)		✓	✓			0.310	0.241	0.527	0.912	3.048
	(2)		✓	✓			0.307	0.239	0.525	0.915	
Team 9	(1)	✓		✓			0.310	0.242	0.515	0.856	3.080
	(2)		✓				0.315	0.239	0.509	0.848	
Reference											3.938
Baseline w/o audio			✓				0.305	0.217	0.481	0.733	
Baseline			✓				0.309	0.215	0.487	0.746	2.848

Table 10: Evaluation results with word-overlapping-based objective measures based on 6 references and a subjective measure based on 5-level ratings for the AVSD track. Under this evaluation, the human rating for the original answers was **3.938**.