

Textual Dataset for Situated Proactive Response Selection

Anonymous ACL submission

Abstract

Recent data-driven conversational models are able to return fluent, consistent, and informative responses to many kinds of requests and utterances in task-oriented situations. However, these responses are typically limited to just the immediate local topic instead of being wider-ranging and proactively taking the conversation further, for example making suggestions proactively to help customers achieve their goals. This inadequacy reflects a lack of understanding of the interlocutor’s situation and implicit goal. To address the problem, we introduce a task of proactive response selection based on situational information. We present a manually-curated dataset of 1.7k English conversation examples that include situational background information plus for each conversation a set of responses, only some of which are acceptable in the situation. A responsive and informed conversation system should select the appropriate responses and avoid inappropriate ones; doing so demonstrates the ability to adequately understand the initiating request and situation. Our benchmark experiments show that this is not an easy task even for strong neural models, offering opportunities for future research. The dataset can be used to develop conversationally informed and proactive dialogue engines. We will release the dataset upon acceptance.

1 Introduction

Conversational assistant systems have recently shown significant improvements for understanding users’ inquiries along with background knowledge, conducting requested operations, and returning natural language responses. Yet, typical systems are likely to be *passive* and only process user-initiated requests or merely ask values for domain-specific slots (Williams et al., 2013; Ammari et al., 2019). In contrast, human assistants like hotel concierges are more *proactive*, acting to address unmentioned needs and expected future events (Cho et al., 1996; Bellini and Convert, 2016). They do not only make

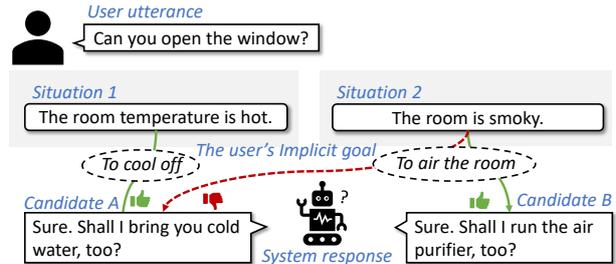


Figure 1: An example of situated goal-aware proactive response selection. The response candidate A is appropriate in Situation 1 but not in Situation 2.

a direct response or a clarification question to their interlocutors but also provide personalized information/assistance based on context and knowledge.

To push the frontier of task-oriented conversation technologies, we propose a task of *proactive* response selection for single-turn help-seeking conversations in English. We mean by proactive that a system engages in an interaction in a cooperative manner (Grice, 1975) and suggests something helpful to a user. The proposed task touches upon two crucial aspects of help-seeking conversations: situation-awareness and goal-awareness.

Situation: Situational information plays an important role in conversations as we illustrate in Figure 1. The example shows a user utterance “Can you open the window?” (top) and two response candidates (bottom). Although both candidates here sound helpful, their appropriateness varies depending on context: When the room is hot, suggesting a cold drink is appropriate assistance (left), but on the other hand, if the room is smoky, then running an air purifier is more helpful (right). Likewise, different situations make different responses more appropriate. A fair amount of situational information can be perceived as visual image, sound, and other kinds of sensory signals, and some of those are effectively incorporated into multi-modal conversational systems (Crook et al., 2019; Kottur et al., 2019). Yet, there are many other types of

information that modern conversation assistance systems have access to, for example, via external APIs such as calendars and maps. In this study, we represent situation statements of six semantic categories (location, possession, etc.) in free English texts, which are more explicit as a semantic representation than just maintaining conversation histories (Lowe et al., 2015; Li et al., 2017; Henderson et al., 2019) and more flexible than structured representations of limited vocabulary (Williams et al., 2013; Budzianowski et al., 2018).

Goal: In the aforementioned example, the two actions address two different goals associated with opening a window, namely, *to cool off* and *to air the room*. While often being unspoken, underlying goals provide important semantic connections among context and utterances on many occasions (Allen and Perrault, 1980) particularly when language is indirect (Perrault, 1980; Walker et al., 2011; Stevens et al., 2015). We use goal information as a stimulus for soliciting naturalistic and proactive responses from human annotators in data collection.

We introduce a dataset of **SitUatated**, **Goal-Aware**, and proactive **Responses** (SUGAR; §3). SUGAR contains 1,761 single-turn English conversation examples, each of which has a user request anchored by an implicit goal, three response candidates with three-point appropriateness ratings, and 12 sentences of situational information. We harvested user utterances and goals from common-sense knowledge bases, ATOMIC (Sap et al., 2019) and ConceptNet (Speer et al., 2017), and collected proactive responses and supporting situational information by crowdsourcing. We then used a language generation model, COMET (Bosselut et al., 2019; Hwang et al., 2021) to generate additional situation statements. Finally, we selected inappropriate (negative) responses for each reference response by an adversarial method to form examples of three-choice response selection. To ensure the data quality, we performed multiple manual validation steps in data collection. Our experiments on SUGAR show that Transformer-based rankers merely achieve Precision@1 of 60-70%, indicating the difficulty of situated reasoning for response selection.

Contributions: (1) We introduce a task of situated proactive response selection and present a high-quality dataset of 1,761 examples by a care-

fully designed data collection pipeline. (2) We perform experiments on our dataset and show that Transformer-based response rankers are not sensitive enough to situational information, which offers opportunities for future research.

2 Related Work

2.1 Conversational Dataset

Acquisition and annotation of real(istic) conversational data has been an essential step for developing conversation engines that imitate human communication (Serban et al., 2018). Various datasets have been constructed with a focus on different aspects of communication.

With regard to target communicative aspects, the most relevant to our work is SIMMC (Moon et al., 2020). SIMMC encompasses surrounding situational information that gives a basis for verbal interactions in task-oriented scenarios in the shopping domain. Moon et al. collected visually-grounded conversation examples from pairs of human annotators interacting with each other in a virtual environment (Crook et al., 2019), where one annotator seeks help for shopping, and the other provides assistance. SUGAR is also concerned with how human interlocutors perform situated conversations in a help-seeking setting. Our work extends this direction to scenarios other than shopping and includes more diverse types of information that modern conversational assistants could access via sensors or external APIs (e.g., temperature and schedule) by representing situational information in a textual form as opposed to visual images.

The choice of modality is motivated by existing conversational datasets that express various kinds of background information in plain text: the persona of an interlocutor (Zhang et al., 2018; Dinan et al., 2020), emotional states (Rashkin et al., 2019), and related documents (Dinan et al., 2019). These examples demonstrate the utility of textual forms for representing both explicit and implicit information of various kinds.

Some existing datasets are concerned with information-seeking conversations like restaurant recommendation where suggestions by assistants naturally occur (e.g., “If you like French cuisine, how about RestaurantX?”, “I can find transportation for you.”). However, it is not trivial to solicit such naturalistic proactive utterances in more diverse help-seeking scenarios. In many cases, the minimum objective of a conversation can be

172 achieved by responding to user-initiated inquiries, 220
173 and such kinds of responses are relatively easy to 221
174 collect from non-expert annotators (Budzianowski 222
175 et al., 2018; Byrne et al., 2019; Eric et al., 2020). 223

176 We address this problem by leveraging implicit 224
177 goals behind user requests. The comprehension 225
178 of goals in conversations has been recognized to 226
179 be important not only in task-oriented dialog re- 227
180 search but also in a broad range of research areas 228
181 such as linguistics, psychology, and artificial in- 229
182 telligence. (Schank and Abelson, 1977; Clark and 230
183 Schaefer, 1989; Gordon and Hobbs, 2004; Rahim- 231
184 toroghi et al., 2017). Human interactions often 232
185 involve indirect speech acts (Perrault, 1980; Gibbs 233
186 and Bryant, 2008) and indirect responses like non- 234
187 yes/no answers to polar questions (Hockey et al., 235
188 1997; de Marneffe et al., 2009; Stevens et al., 2015; 236
189 Louis et al., 2020). These studies motivate our 237
190 strategy for soliciting natural-sounding proactive 238
191 responses from crowd workers. 239

192 In contrast to most datasets we introduced here, 240
193 SUGAR only contains single-turn conversation ex- 241
194 amples due to the ease of data collection and quality 242
195 control. Yet, our primary focus is on conversational 243
196 assistance, which is likely to be completed within 244
197 a few turns (Völkel et al., 2021). Thus, we believe 245
198 that single-turn examples are still useful for system 246
199 development. It is possible to extend our problem 247
200 setting and data collection approach to a multi-tern 248
201 setting, which we leave as future work. 249

202 2.2 Response Selection 250

203 Automatic response models can be divided into 251
204 two approaches: response generation and response 252
205 selection. Response generation directly generates 253
206 natural language response text from scratch, and 254
207 response selection selects a response from a candi- 255
208 date pool built by humans, templates, or language 256
209 generation systems. The latter approach is widely 257
210 used in many real-world applications cases because 258
211 of the controllability of responses and the easiness 259
212 of evaluation (Deriu et al., 2020). In this study, we 260
213 study the task of response selection. We assume 261
214 that an external response generation system gener- 262
215 ates candidates based on the system’s functionality 263
216 and focus on picking the appropriate ones.¹

217 To train and evaluate a response selection system, 264
218 each example must have distractors (negative re- 265
219 sponses), but typically, conversational datasets only

¹Training and evaluating generation systems on SUGAR is also possible and is an interesting direction for future research.

220 contain ground truth responses. Thus, it has been 221
222 commonly practiced to pick negative responses 223
224 by random sampling (Lowe et al., 2015; Hender- 225
226 son et al., 2019). This approach is practically ad- 227
228 vantageous but may introduce negative responses 229
230 that are clearly off-topic or false negatives (Akama 230
231 et al., 2020; Hedayatnia et al., 2022). To allevi- 231
232 ate this problem, we use an adversarial filtering 232
233 algorithm (Zellers et al., 2018; Sakaguchi et al., 233
234 2019; Bhagavatula et al., 2020) to select compet- 234
235 itive distractors and recruit crowd workers to rate 235
236 candidates, allowing each example to have multiple 236
237 acceptable responses. 237

238 3 Task and Data 239

238 The goal of this study is to provide a resource for 239
239 developing a system that can observe situational 240
240 information and return a proactive response to a 241
241 user. We consider six categories of observable *sit-* 242
242 *uational information*²: location (where the user 243
243 is), possession (what the user has), time, date, be- 244
244 havior (what the user is/was doing), and environ- 245
245 ment (temperature, etc.) We define a *proactive* 246
246 response to be a response that provides *suggestions* 247
247 to help users achieve their goals. 248

249 3.1 Problem Formulation 250

249 Our task has five components: (1) a user utterance 251
250 u , (2) situation statements $S = \{s_i\}_{i=1,\dots,l}$, where 252
251 l is the number of statements, (3) responses $R =$ 253
252 $\{r_i\}_{i=1,\dots,m}$, where m is the number of response 254
253 candidates³, (4) their appropriateness ratings $Y =$ 255
254 $\{y_i\}_{i=1,\dots,m}$, where y_i is a three-point Likert scale, 256
255 and (5) an implicit goal g . S can include distractors 257
256 that are not directly relevant to the conversation. u , 258
257 S , and R are given as input, and the task is to re- 259
258 rank R . Response selection systems are trained and 260
259 evaluated by Y . 261

262 3.2 Data 263

262 SUGAR contains 1,761 high-quality examples, 264
263 each of which has three response candidates and 12 265
264 sentences of situational information. Table 1 shows 266
265 the dataset statistics. We constructed the dataset 267
266 with the eight steps shown in Figure 2. We describe 268
267 them below.⁴ 269

²See Appendix D for more details.

³We pick $m - 1$ responses automatically such that they are less appropriate than the reference response in a given context (See Step 7). Nevertheless, there usually exist one or more acceptable responses to a given user utterance. We thus annotate all acceptable responses manually (Step 8).

⁴See also Appendix B for more details.

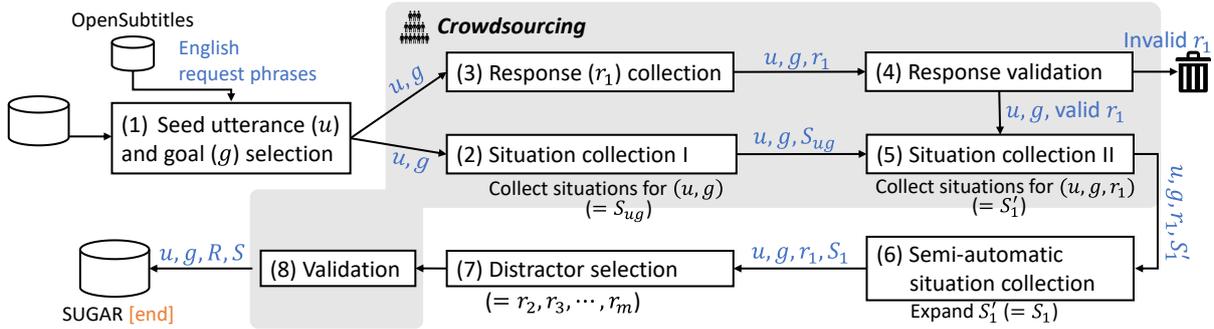


Figure 2: Pipeline for data collection. We start with existing common-sense knowledge bases (ATOMIC and ConceptNet) and extract utterance and goal events as a seed (1). We collect responses and situation statements for each seed by crowdsourcing (2-5), acquire more situation statements semi-automatically (6), and select distractor responses and situations to form response selection examples (7). We finally validate the examples manually (8). Steps (2) and (3-4) are executed in parallel.

	u	r	g	s
Sent types.	370	1,754	431	4,363
Tokens	12,172	29,389	7,501	146,526
Avg. tokens per ex.	6.9	16.7	4.3	83.2

Table 1: Dataset statistics. The total number of examples is 1,761 (34,590 sentences).

(1) Seed Utterance & Goal Selection: We harvested action and goal events from two common-sense knowledge bases, ATOMIC (Sap et al., 2019) and ConceptNet (Speer et al., 2017), where knowledge is represented as nodes representing events or concepts and edges connecting them with semantic relations. The collected action-goal node pairs served as the seed utterance-goal for soliciting responses and situational statements in the following data collection steps. First, we extracted nodes consisting of verb phrases (VPs) that appear at least five times within English request phrases (e.g., Please VP, Could you VP?, etc.) in the OpenSubtitles corpus (Henderson et al., 2019). These request expressions were also used as the surface form of u . Two of the authors then selected 563 events that can be achieved within a reasonable time span, can be assisted by someone else, and can be triggered by a goal. We retrieved their implicit goals g by goal-related edges in ATOMIC and ConceptNet. Specifically, we used `xNeed` in the reverse direction and `xIntent` in ATOMIC and `HasPrerequisite` in the reverse direction and `MotivatedByGoal` in ConceptNet. Finally, two of the authors evaluated the node pairs and picked 501 (u, g) pairs for which we can naturally say “I do u to achieve g .” We also merged synonymous expressions (e.g., *go to a market* and *go to a supermarket*) into a

single entry and corrected grammatical errors and unnatural phrases.

(2) Situation Collection I: We collected situation statements in two phases to simplify annotation work. The first phase focuses on u and g , and the second phase considers r in addition to u and g . In this step, we presented a pair of u and g texts to crowd workers and instructed them to specify situational information that is required to guess the goal based on the utterance. For example, an implicit goal “to cool off” can be naturally inferred by situations like “The user is home. The room temperature is hot.” We asked workers to write *observable* facts in the six semantic categories (§3). For example, “The room temperature is hot.” is valid, but “The user feels hot.” is invalid as assistance systems cannot *observe* the user’s feeling. We recruited one worker for each (u, g) pair and paid \$0.12 per HIT⁵ (one (u, g) pair/HIT).

(3) Response Collection: In parallel to Step (2), we recruited two crowd workers for each (u, g) pair to collect responses. The workers created at least two responses: one of the responses accepts and the other rejects the request. We asked the workers to write a *proactive* response, a response providing suggestions for goal fulfilment.⁶ To solicit responses closely connected to implicit goals rather than to domain knowledge, we instructed the workers to avoid posing a clarification question like “Sure, I’ll turn on the air conditioner for you.

⁵Human Intelligence Task, a unit of task in MTurk.

⁶For a response that rejects a user’s request, we instructed the workers to provide a reason for rejection (e.g., we cannot brew coffee *because we are out of coffee filters*) in addition to a suggestion.

Loc.	Poss.	Time	Date	Behav.	Env.
1,991	3,567	1,084	149	1,688	2,794

Table 2: Number of situation statements ($\in S_1$).

321 *Would you like it on a high or low setting?* (= clarifi-
322 cation) The workers were presented one u - g pair
323 in each HIT and were paid \$0.30/HIT.

324 **(4) Response Validation:** We present the utterances,
325 goals, and collected responses to crowd workers and
326 evaluated the helpfulness of the response. A response
327 is considered to be valid if it satisfies the following
328 criteria: (1) the response suggests or requests some-
329 thing new, and (2) the suggestion or request is help-
330 ful for achieving the goal. Each response was evalu-
331 ated by three workers. We then picked the responses
332 that were approved by two or three workers. We call
333 a verified response a *reference response* r_1 hereaf-
334 ter. Each HIT contains up to seven responses, and one
335 of them is a dummy question for evaluating crowd
336 workers. For quality control, we filtered out crowd
337 workers who participated in the task twice or more
338 and did not reach 0.75% accuracy for the dummy
339 questions. The workers were paid \$0.18 for this task.

341 **(5) Situation Collection II:** We collected situa-
342 tion statements from crowd workers with the follow-
343 ing two goals: (1) to collect situation statements
344 that cover the reference response r_1 and (2) to ver-
345 ify the situation statements collected in Step (2).
346 We presented (u, g, r_1) with the statements ob-
347 tained in Step (2) and again instructed crowd work-
348 ers to write observable facts. The results of Step
349 (2) were provided as editable initial values, and we
350 encouraged workers to update the texts when it is
351 necessary. We recruited one crowd worker for each
352 (u, g, r_1) with the reward of \$0.42/HIT.

353 **(6) Semi-automatic Situation Collection:** We
354 found that the collected situations were often under-
355 or over-specified. We addressed this by automatic
356 situation generation and manual verification.

357 One of the authors examined all the situation
358 statements, discarded/modified inappropriate situa-
359 tions, and categorized them into six categories. We
360 then used the cleaned and labeled texts to fine-tune
361 a neural sequence-to-sequence to generate more sit-
362 uations. Specifically, we fine-tuned BART (Lewis
363 et al., 2020) trained on ATOMIC₂₀ (Hwang et al.,

Input

[u]Please open the window. (u text)

[g]to cool off (g text)

[r]Sure, shall I bring cold water, too? (r_1 text)

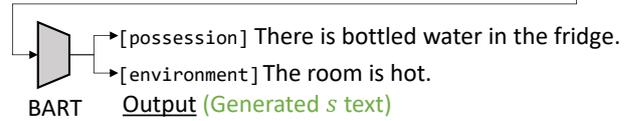


Figure 3: Example of automatic situation generation by BART (Step 6). [$*$] is a special symbol to denote the type of the following text. The first token of output is given as a prompt to control the semantic category of output.

2021)⁷ to take a concatenation of u , g , and r_1 as
364 input and generate a text for a given situation cat-
365 egory as illustrated in Figure 3. We performed a
366 beam search of width 3 and took top-3 generation
367 results for each input and relation. Finally, we man-
368 ually verified the generated situations, resulting
369 in 4,375 unique situations (6.4 ± 1.3 statements
370 per example). We denote the situation statements
371 attached to (u, g, r_1) by S_1 . Table 2 shows the dis-
372 tribution of situation categories in SUGAR. State-
373 ments about possession and environment appear
374 most frequently, which is reasonable because such
375 situational information often decides actions that
376 can be carried out (e.g., to drink coffee, coffee
377 must be available). The other categories are less
378 frequent, but 64% of examples have at least one
379 time or date information, and 69% have a statement
380 about behavior.

381 **(7) Distractor Selection:** The examples col-
382 lected in the previous steps only contain reference
383 responses r_1 and supporting situation statements
384 S_1 . We added $m - 1$ response candidates along
385 with their relevant situational information as dis-
386 tractors so that all examples have m response
387 candidates and l situation statements. We set $m = 3$
388 and $l = 12$. In this section, we describe the high-
389 level idea of our algorithm. Appendix C presents
390 technical details.

391 Distractors can be obtained by random sampling
392 as practiced in many studies (Henderson et al.,
393 2019) or by advanced methods such as adversarial
394 filtering (Li et al., 2019; Gupta et al., 2021).
395 However, such approaches may introduce off-topic
396 responses that are easy to rule out and false nega-
397

⁷Note that the framework of pre-training Transformer models on common-sense knowledge bases was originally proposed by Bosselut et al. (2019).

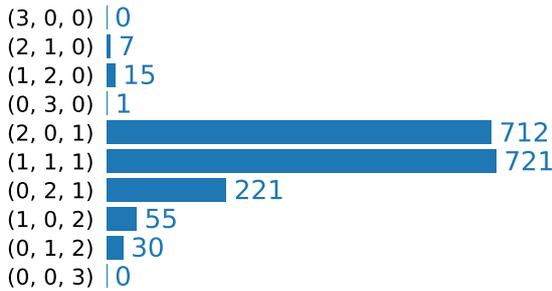


Figure 4: Result of rating annotations (Step 8). The labels denote (the number of *Bad* options, the number of *Acceptable* options, number of *Best* options). We removed one example with three *Acceptable* responses (0, 3, 0) from the dataset.

tives — acceptable responses treated as negative examples, degrading system performance as well as reliability of evaluation (Akama et al., 2020; Hedayatnia et al., 2022).

To alleviate this problem, we combine lexical matching and adversarial filtering (Zellers et al., 2018; Sakaguchi et al., 2019; Bhagavatula et al., 2020) to construct distractors and validate them manually (see Step 8). We first created an initial dataset by a lightweight method based on sentence embeddings and lexical matching. We then performed $J = 3$ rounds of adversarial filtering. In each round, we split the dataset into $K = 10$ folds, and for each split, we trained a binary logistic regression classifier that takes sentence embeddings of u , S_1 , and a response candidate. We computed sentence embeddings by SentenceTransformers (Reimers and Gurevych, 2019) with MP-Net (Song et al., 2020). We used the trained classifier to identify easy distractors and replace them with more confusing ones with respect to the score function. We sampled $m - 1 = 2$ responses for each example. All response candidates in the same example have the same polarity. Finally, we expanded S_1 , which only contains relevant information to u and r_1 , to obtain a set of $l = 12$ situations S such that some of them are related to distractors but do not disqualify r_1 , and statements do not contradict with each other. We again used sentence embeddings to find topically related situational information and avoid contradiction with keyword-based heuristics.

(8) Validation: There are usually multiple appropriate responses in one conversational context, and therefore, some of the challenging “distractors” picked in the previous step can be acceptable

or even more appropriate than the reference r_1 . To avoid introducing false negatives, we rated all response candidates on a three-point Likert scale (*Bad*, *Acceptable*, or *Best*) by crowdsourcing. We recruited three crowd workers per example with the reward of \$0.25/each and asked them to pick an appropriate response candidate (Krippendorff’s α (Krippendorff, 2004) of 0.484). We then aggregated ratings by the statistical model proposed by Zhou et al. (2014) to obtain the final rating Y .⁸ We discarded one example in this validation step and obtained 1861 examples with all responses rated. Figure 4 shows the annotation result. As we expected, a fair number of examples (56%) have more than one *Best* or *Acceptable* responses.

4 Experiments

We evaluate several baseline models on SUGAR to explore two questions concerned with the nature of the proposed task and dataset: (1) Is understanding of situational information required to identify proactive responses in SUGAR? (2) Can standard matching-based systems capture relevant situational information and solve the task?

4.1 Baselines

We evaluate a lexical-matching approach and several Transformer-based response selection systems. A variety of neural networks have been proposed for the task of response selection Tao et al. (2021), but we opted to focus on the direct application of pre-trained Transformers rather than equipping them with extra modules/resources as Transformers have proven effective in conversation tasks with minimal adaptation (Budzianowski and Vulić, 2019; Han et al., 2021) and are also used as a basis of many recent complex approaches.

TF-IDF ranker: As the simplest approach in our evaluation, we used a lexical-matching baseline system that ranks response candidates by cosine similarity of TF-IDF vectors of context and a response candidate (Lowe et al., 2015). We calculated TF-IDF weights on a training split with `scikit-learn` library.

Transformer ranker: We fine-tuned and evaluated three variants of Transformer-based rankers:

⁸In the first run, all candidates were rated as equally good or bad in 18 examples. We updated and re-annotated 17 examples.

- 479 1. **BERT-FP:** We fine-tuned the BERT-FP
 480 model (Han et al., 2021) trained on Ubuntu Di-
 481 alogue Corpus (Lowe et al., 2015) on SUGAR.
 482 BERT-FP is a response selection model based
 483 on uncased BERT_{base} (Devlin et al., 2019), a
 484 12-layer Transformer model ($\approx 110\text{M}$ param-
 485 eters) pre-trained on BookCorpus and English
 486 Wikipedia. Han et al. fine-tuned the model on
 487 Ubuntu Dialogue Corpus (Lowe et al., 2015)
 488 by fine-grained post-training and supervised
 489 learning and reported high performance on the
 490 Ubuntu dataset.
- 491 2. **BERT:** We also fine-tuned BERT_{base} without
 492 the additional training mentioned above in or-
 493 der to analyze the benefit of the additional
 494 training on conversational texts in a different
 495 domain. We used a cased model following
 496 BERT-FP. In the experiments of Hedayatnia
 497 et al. (2022), the BERT ranker performed on
 498 par with BERT-FP.
- 499 3. **RoBERTa:** To analyze the effect of an un-
 500 derlying Transformer model, we evaluated
 501 cased RoBERTa_{base} (Liu et al., 2019) to rank
 502 response candidates in the same way we did
 503 with BERT. RoBERTa_{base} is a 12-layer Trans-
 504 former model ($\approx 125\text{M}$ parameters) trained
 505 with an improved training configuration, and
 506 Liu et al. report that RoBERTa outperforms
 507 BERT in multiple tasks and datasets.

508 We optimize model parameters with
 509 Adam (Kingma and Ba, 2015) to minimize
 510 max-margin loss.

511 4.2 Experimental Setup

512 **Input format:** We concatenated context and a
 513 response candidate for the Transformer rankers. To
 514 address our questions, we experimented with three
 515 variants of context:

- 516 1. u : Utterance (u)-only
 517 2. $u + S_1$: Utterance (u) plus relevant situation
 518 (S_1)
 519 3. $u + S$: Utterance (u) plus relevant and irrele-
 520 vant situation (S)

521 **Training and Test:** We performed five-fold
 522 cross-validation (training:validation:test=6:2:2).⁹
 523 For each round, we trained a Transformer ranker
 524 for 10 epochs with a batch size of 32 and evaluated

⁹We removed examples without *Bad* response options from the validation and test splits

System	Context	Prec.@1	nDCG@3
TF-IDF	u	.6044 \pm .0390	.8329 \pm .0136
	$u + S_1$.7541 \pm .0106	.9130 \pm .0037
	$u + S$.5350 \pm .0372	.8408 \pm .0124
BERT-FP	u	.6269 \pm .0184	.8732 \pm .0078
	$u + S_1$.7728 \pm .0137	.9262 \pm .0027
	$u + S$.6371 \pm .0219	.8764 \pm .0070
BERT	u	.7053 \pm .0099	.9007 \pm .0054
	$u + S_1$.8128 \pm .0135	.9408 \pm .0056
	$u + S$.7081 \pm .0232	.9044 \pm .0081
RoBERTa	u	.7155 \pm .0113	.9059 \pm .0040
	$u + S_1$.8140 \pm .0096	.9432 \pm .0026
	$u + S$.7047 \pm .0152	.9065 \pm .0056

Table 3: Average test scores over five-fold cross-validation. Prec.@1 stands for Top-1 precision.

525 the model by nDCG@3 on the validation split ev-
 526 ery epoch. We then selected the best checkpoint for
 527 evaluation. To stabilize training, we applied weight
 528 decay of 0.05, set the maximum gradient norm to
 529 5.0, and used a linear learning rate scheduler with
 530 5% (≈ 20) warm-up steps. We further performed
 531 light-weight grid-search for hyperparameter tuning
 532 based on an average nDCG@3 score on validation
 533 splits, with learning rate $\in \{5e - 5, 1e - 5\}$, and
 534 margin for the max-margin loss $\in \{1.0, 0.5, 0.1\}$.
 535 One epoch of training took 1-2m on GeForce GTX
 536 TITAN X. We report the average Precision@1 and
 537 nDCG@3 on the test splits.

538 4.3 Results

539 Table 3 shows the average test scores over five-fold
 540 cross-validation. We can see two general patterns:
 541 (1) The transformer-based models except for BERT-
 542 FP outperformed the TF-IDF baseline, and (2) the
 543 systems that are only given relevant statements S_1
 544 along with u outperformed their counterparts with
 545 different context settings. With regard to the key
 546 questions, the result provides several interesting
 547 findings:

- 548 1. Comparing two context settings u and $u + S_1$,
 549 we can see that relevant situational informa-
 550 tion brings in a clear performance boost as
 551 we expected (e.g., +0.11 in Precision@1 and
 552 +0.04 in nDCG@3 with BERT).
 553 2. The performance gain from S_1 can be at-
 554 tributed to increased word overlaps between
 555 context and the right responses as the result

of TF-IDF indicates. However, with a few distractor statements added in the $u + S$ setting, the performance of the TF-IDF baseline dropped substantially (-0.22 in Precision@1 and -0.07 in nDCG@3). This means that our dataset successfully avoids superficial clues, calling for a higher-level understanding of situational information.

3. Interestingly, the performance of Transformer rankers also decreased drastically, to the same level as their corresponding systems without situational statements in input.
4. The additional pre-training of BERT-FP was not beneficial to our task, which is consistent with the observation of Hedayatnia et al. (2022). We speculate that this is due to the domain mismatch of training corpora. BERT-FP is pre-trained on technical topics related to Ubuntu, but SUGAR is concerned with a wider range of topics in daily life.

These observations provide insights into our questions. First, the understanding of situational information is necessary for accurately selecting proactive responses, indicating that SUGAR is an effective resource to develop and evaluate situated conversation systems. Second, it is not trivial for Transformer rankers to pick up useful clues out of a mixture of relevant and irrelevant situation statements.

5 Limitations

Data size: SUGAR is relatively small compared to recently published datasets. This is due to the complexity of our problem setting and annotation pipeline. We prioritized quality over quantity and performed multiple steps of manual intervention to reduce errors, false negatives, and annotation artifacts. These problems have been reported in various NLP tasks not limited to conversational tasks (Gururangan et al., 2018; Akama et al., 2020; Elazar et al., 2020). Nonetheless, our experiment has shown that pre-trained Transformer models can be trained to outperform a TF-IDF ranker by a clear margin, which is encouraging. In addition, we could automatically induce noisy but large-scale training instances from existing resources, for example, by harvesting event pairs that can be used as u and r from event knowledge bases such as ATOMIC₂₀ and generating situation statements using our generator (§3).

Representation of situation information: In SUGAR, situation information is represented in textual expressions. In real-world applications, such information could be collected via external APIs (e.g., calendar and map) and sensors (e.g., camera) and stored in non-textual forms. Our study is a proof-of-concept that shows the understanding of situational information is very important for response selection. Future research should explore ways to process situation information that is expressed in other forms of data (e.g., structured texts, numbers, images). Even if the value is structured or images, we could transform them into textual forms as done in data-to-text research (Shen et al., 2020; Miura et al., 2021). Besides, we acknowledge that situational information is often under-specified in SUGAR because some information is considered to be common-sense (e.g., a room has a door) or presupposed (e.g., “Please open the door” presupposes that the door is closed.), and such information was not explicitly stated by human annotators during data collection. Therefore, response selection systems should be equipped with a mechanism to handle implicit knowledge to solve the task.

6 Conclusion and Future Work

We proposed a task of situated proactive response selection for developing and evaluating conversational assistants that can help users proactively in various help-seeking scenarios. We constructed a dataset of 1.7k examples by crowdsourcing and semi-automatic generation.

There are several interesting directions for future research. First, as shown in our experiments, it is challenging to pick up relevant situational information and use it to reason about user requests and potential assistance. To achieve this, conversational systems will need to be equipped with world knowledge to effectively align situation information with an interaction. One promising approach is knowledge-based response models such as graph neural networks, which recently has shown to be effective in various NLP tasks (Zhang et al., 2020; Zhou et al., 2022; *inter alia*). Second, although we leveraged implicit goals only for soliciting proactive responses in data collection in this study, understanding of goals should be necessary for building better conversation engines as claimed in early studies (Allen and Perrault, 1980; *inter alia*). We believe SUGAR can facilitate future research in this direction.

655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709

References

Reina Akama, Sho Yokoi, Jun Suzuki, and Kentaro Inui. 2020. [Filtering noisy dialogue corpora by connectivity and content relatedness](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing*, pages 941–958. Online. Association for Computational Linguistics.

James F. Allen and C. Raymond Perrault. 1980. [Analyzing intention in utterances](#). *Artificial Intelligence*, 15(3):143–178.

Tawfiq Ammari, Jofish Kaye, Janice Y. Tsai, and Frank Bentley. 2019. [Music, search, and IoT: How people \(really\) use voice assistants](#). *ACM Transactions on Computer-Human Interaction*, 26(3):17:1–17:28.

Nicola Bellini and Laetitia Convert. 2016. [The concierge. tradition, obsolescence and innovation in tourism](#). *Symphonya. Emerging Issues in Management*, 0(2):17.

Chandra Bhagavatula, Ronan Le Bras, Chaitanya Malaviya, Keisuke Sakaguchi, Ari Holtzman, Hannah Rashkin, Doug Downey, Scott Wen-tau Yih, and Yejin Choi. 2020. [Abductive commonsense reasoning](#).

Antoine Bosselut, Hannah Rashkin, Maarten Sap, Chaitanya Malaviya, Asli Celikyilmaz, and Yejin Choi. 2019. [COMET: Commonsense Transformers for automatic knowledge graph construction](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4762–4779, Florence, Italy. ArXiv: 1906.05317.

Paweł Budzianowski and Ivan Vulić. 2019. [Hello, it’s GPT-2 - how can i help you? towards the use of pre-trained language models for task-oriented dialogue systems](#). In *Proceedings of the 3rd Workshop on Neural Generation and Translation*, pages 15–22, Hong Kong. Association for Computational Linguistics.

Paweł Budzianowski, Tsung-Hsien Wen, Bo-Hsiang Tseng, Iñigo Casanueva, Stefan Ultes, Osman Ramadan, and Milica Gašić. 2018. [MultiWOZ - A large-scale multi-domain Wizard-of-Oz dataset for task-oriented dialogue modelling](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 5016–5026, Brussels, Belgium. Association for Computational Linguistics.

Bill Byrne, Karthik Krishnamoorthi, Chinnadhurai Sankar, Arvind Neelakantan, Ben Goodrich, Daniel Duckworth, Semih Yavuz, Amit Dubey, Kyu-Young Kim, and Andy Cedilnik. 2019. [Taskmaster-1: Toward a realistic and diverse dialog dataset](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing*, pages 4516–4525, Hong Kong, China. Association for Computational Linguistics.

Wonae Cho, Robert T Sumichrast, and Michael D Olsen. 1996. [Expert-system technology for hotels: Concierge application](#). *Cornell Hotel and Restaurant Administration Quarterly*, 37(1):54–60. 710
711
712
713

Herbert H. Clark and Edward F. Schaefer. 1989. [Collaborating on contributions to conversations](#). In *Language Processing in Social Context*, volume 54 of *North-Holland Linguistic Series: Linguistic Variations*, pages 123–152. Elsevier. ISSN: 0078-1592. 714
715
716
717
718

Paul A. Crook, Shivani Poddar, Ankita De, Semir Shafi, David Whitney, Alborz Geramifard, and Rajen Subba. 2019. [SIMMC: Situated interactive multi-modal conversational data collection and evaluation platform](#). *IEEE Workshop on Automatic Speech Recognition and Understanding*. 719
720
721
722
723
724

Marie-Catherine de Marneffe, Scott Grimm, and Christopher Potts. 2009. [Not a simple yes or no: Uncertainty in indirect answers](#). In *Proceedings of the SIGDIAL 2009 Conference*, pages 136–143, London, UK. Association for Computational Linguistics. 725
726
727
728
729

Jan Deriu, Alvaro Rodrigo, Arantxa Otegi, Guillermo Echevoyen, Sophie Rosset, Eneko Agirre, and Mark Cieliebak. 2020. [Survey on evaluation methods for dialogue systems](#). *Artificial Intelligence Review*, pages 1–56. 730
731
732
733
734

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional Transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics. 735
736
737
738
739
740
741
742

Emily Dinan, Varvara Logacheva, Valentin Malykh, Alexander Miller, Kurt Shuster, Jack Urbanek, Douwe Kiela, Arthur Szlam, Iulian Serban, Ryan Lowe, Shrimai Prabhumoye, Alan W. Black, Alexander Rudnicky, Jason Williams, Joelle Pineau, Mikhail Burtsev, and Jason Weston. 2020. [The second conversational intelligence challenge \(ConvAI2\)](#). In *The NeurIPS ’18 Competition*, The Springer Series on Challenges in Machine Learning, pages 187–208, Cham. Springer International Publishing. 743
744
745
746
747
748
749
750
751
752

Emily Dinan, Stephen Roller, Kurt Shuster, Angela Fan, Michael Auli, and Jason Weston. 2019. [Wizard of Wikipedia: Knowledge-powered conversational agents](#). In *The Seventh International Conference on Learning Representations*, New Orleans, Louisiana, USA. 753
754
755
756
757
758

Yanai Elazar, Victoria Basmov, Shauli Ravfogel, Yoav Goldberg, and Reut Tsarfaty. 2020. [The extraordinary failure of complement coercion crowdsourcing](#). In *Proceedings of the First Workshop on Insights from Negative Results in NLP*, pages 106–116, Online. Association for Computational Linguistics. 759
760
761
762
763
764

765	Mihail Eric, Rahul Goel, Shachi Paul, Abhishek Sethi,	Beth Ann Hockey, Deborah Rossen-Knill, Beverly Spe-	820
766	Sanchit Agarwal, Shuyang Gao, Adarsh Kumar, Anuj	jewski, Matthew Stone, and Stephen Isard. 1997. Can	821
767	Goyal, Peter Ku, and Dilek Hakkani-Tur. 2020. Mul-	you predict answers to yes/no questions? Yes, no and	822
768	tiWOZ 2.1: A consolidated multi-domain dialogue	stuff. In <i>Proceedings of the Fifth European Con-</i>	823
769	dataset with state corrections and state tracking base-	<i>ference on Speech Community and Technology (EU-</i>	824
770	lines . In <i>Proceedings of The 12th Language Re-</i>	<i>ROSPEECH)</i> , pages 2267–2270, Rhodes, Greece.	825
771	<i>sources and Evaluation Conference</i> , pages 422–428,		
772	Marseille, France. European Language Resources	Jena D. Hwang, Chandra Bhagavatula, Ronan Le Bras,	826
773	Association.	Jeff Da, Keisuke Sakaguchi, Antoine Bosselut, and	827
		Yejin Choi. 2021. COMET-ATOMIC 2020: On sym-	828
774	Raymond W Gibbs and Gregory A Bryant. 2008. Striv-	bolic and neural commonsense knowledge graphs . In	829
775	ing for optimal relevance when answering questions .	<i>Proceedings of the Thirty-Fifth AAAI Conference on</i>	830
776	<i>Cognition</i> , 106(1):345–369.	<i>Artificial Intelligence</i> , pages 6384–6392, Online.	831
777	Andrew S. Gordon and Jerry R. Hobbs. 2004. Formal-	Diederik P. Kingma and Jimmy Ba. 2015. Adam: A	832
778	izations of commonsense psychology . <i>AI Magazine</i> ,	method for stochastic optimization . <i>The Third Inter-</i>	833
779	25(4):49.	<i>national Conference for Learning Representations</i> .	834
780	H. P. Grice. 1975. Logic and conversation . In <i>Speech</i>	Satwik Kottur, José M. F. Moura, Devi Parikh, Dhruv	835
781	<i>Acts</i> , number 3 in Syntax and Semantics, pages 41 –	Batra, and Marcus Rohrbach. 2019. CLEVR-Dialog:	836
782	58. Academic Press, New York, NY.	A diagnostic dataset for multi-round reasoning in vi-	837
		sual dialog . In <i>Proceedings of the 2019 Conference</i>	838
783	Prakhar Gupta, Yulia Tsvetkov, and Jeffrey Bigham.	<i>of the North American Chapter of the Association for</i>	839
784	2021. Synthesizing adversarial negative responses	<i>Computational Linguistics: Human Language Tech-</i>	840
785	for robust response ranking and evaluation . In <i>Find-</i>	<i>nologies</i> , pages 582–595, Minneapolis, Minnesota.	841
786	<i>ings of the Association for Computational Linguis-</i>	Association for Computational Linguistics.	842
787	<i>tics: ACL-IJCNLP 2021</i> , pages 3867–3883, Online.		
788	Association for Computational Linguistics.	Klaus Krippendorff. 2004. Measuring the reliability of	843
		qualitative text analysis data . <i>Quality & Quantity</i> ,	844
789	Suchin Gururangan, Swabha Swayamdipta, Omer Levy,	38(6):787–800.	845
790	Roy Schwartz, Samuel Bowman, and Noah A. Smith.		
791	2018. Annotation artifacts in natural language infer-	Mike Lewis, Yinhan Liu, Naman Goyal, Marjan	846
792	ence data . In <i>Proceedings of the 2018 Conference of</i>	Ghazvininejad, Abdelrahman Mohamed, Omer Levy,	847
793	<i>the North American Chapter of the Association for</i>	Veselin Stoyanov, and Luke Zettlemoyer. 2020.	848
794	<i>Computational Linguistics: Human Language Tech-</i>	BART: Denoising sequence-to-sequence pre-training	849
795	<i>nologies</i> , pages 107–112, New Orleans, Louisiana.	for natural language generation, translation, and com-	850
796	Association for Computational Linguistics.	prehension . In <i>Proceedings of the 58th Annual Meet-</i>	851
		<i>ing of the Association for Computational Linguistics</i> ,	852
797	Janghoon Han, Taesuk Hong, Byoungjae Kim,	pages 7871–7880, Online. Association for Computa-	853
798	Youngjoong Ko, and Jungyun Seo. 2021. Fine-	tional Linguistics.	854
799	grained post-training for improving retrieval-based		
800	dialogue systems . In <i>Proceedings of the 2021 Con-</i>	Jia Li, Chongyang Tao, Wei Wu, Yansong Feng,	855
801	<i>ference of the North American Chapter of the Asso-</i>	Dongyan Zhao, and Rui Yan. 2019. Sampling mat-	856
802	<i>ciation for Computational Linguistics: Human Lan-</i>	ters! an empirical study of negative sampling strate-	857
803	<i>guage Technologies</i> , pages 1549–1558, Online. As-	gies for learning of matching models in retrieval-	858
804	association for Computational Linguistics.	based dialogue systems . In <i>Proceedings of the 2019</i>	859
		<i>Conference on Empirical Methods in Natural Lan-</i>	860
805	Behnam Hedayatnia, Di Jin, Yang Liu, and Dilek	<i>guage Processing and the 9th International Joint</i>	861
806	Hakkani-Tur. 2022. A systematic evaluation of re-	<i>Conference on Natural Language Processing</i> , pages	862
807	sponse selection for open domain dialogue . In <i>Pro-</i>	1291–1296, Hong Kong, China. Association for Com-	863
808	<i>ceedings of the 23rd Annual Meeting of the Special</i>	putational Linguistics.	864
809	<i>Interest Group on Discourse and Dialogue</i> , pages		
810	298–311, Edinburgh, UK. Association for Computa-	Yanran Li, Hui Su, Xiaoyu Shen, Wenjie Li, Ziqiang	865
811	tional Linguistics.	Cao, and Shuzi Niu. 2017. DailyDialog: A manually	866
		labelled multi-turn dialogue dataset . In <i>Proceedings</i>	867
812	Matthew Henderson, Iñigo Budzianowski	<i>of the Eighth International Joint Conference on Nat-</i>	868
813	Pawełand Casanueva, Sam Coope, Daniela	<i>ural Language Processing</i> , pages 986–995, Taipei,	869
814	Gerz, Girish Kumar, Nikola Mrkšić, Georgios Sp-	Taiwan. Asian Federation of Natural Language Pro-	870
815	ithourakis, Pei-Hao Su, Ivan Vulčić, and Tsung-Hsien	cessing.	871
816	Wen. 2019. A repository of conversational datasets .		
817	In <i>Proceedings of the First Workshop on NLP for</i>	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Man-	872
818	<i>Conversational AI</i> , pages 1–10, Florence, Italy.	dar Joshi, Danqi Chen, Omer Levy, Mike Lewis,	873
819	Association for Computational Linguistics.	Luke Zettlemoyer, and Veselin Stoyanov. 2019.	874
		RoBERTa: A robustly optimized BERT pretraining	875
		approach . <i>arXiv</i> .	876

877	Annie Louis, Dan Roth, and Filip Radlinski. 2020. “I’d rather just go to bed”: Understanding indirect answers. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing</i> , pages 7411–7425, Online. Association for Computational Linguistics.	934
878		935
879		936
880		937
881		
882		
883	Ryan Lowe, Nissan Pow, Iulian Serban, and Joelle Pineau. 2015. The Ubuntu dialogue corpus: A large dataset for research in unstructured multi-turn dialogue systems. In <i>Proceedings of the 16th Annual Meeting of the Special Interest Group on Discourse and Dialogue</i> , pages 285–294, Prague, Czech Republic. Association for Computational Linguistics.	938
884		939
885		940
886		941
887		942
888		943
889		944
890	Yasuhide Miura, Yuhao Zhang, Emily Tsai, Curtis Langlotz, and Dan Jurafsky. 2021. Improving factual completeness and consistency of image-to-text radiology report generation. In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 5288–5304, Online. Association for Computational Linguistics.	945
891		
892		
893		
894		
895		
896		
897		
898	Seungwhan Moon, Satwik Kottur, Paul Crook, Ankita De, Shivani Poddar, Theodore Levin, David Whitney, Daniel Difranco, Ahmad Beirami, Eunjoon Cho, Rajen Subba, and Alborz Geramifard. 2020. Situated and interactive multimodal conversations. In <i>Proceedings of the 28th International Conference on Computational Linguistics</i> , pages 1103–1121, Barcelona, Spain (Online). International Committee on Computational Linguistics.	946
899		947
900		948
901		949
902		
903		
904		
905		
906		
907	C. Raymond Perrault. 1980. A plan-based analysis of indirect speech act. <i>American Journal of Computational Linguistics</i> , 6(3-4):167–182.	950
908		951
909		952
910	Elahe Rahimtoroghi, Jiaqi Wu, Ruimin Wang, Pranav Anand, and Marilyn Walker. 2017. Modelling protagonist goals and desires in first-person narrative. In <i>Proceedings of the 18th Annual SIGdial Meeting on Discourse and Dialogue</i> , pages 360–369, Saarbrücken, Germany. Association for Computational Linguistics.	953
911		954
912		955
913		956
914		957
915		958
916		959
917	Hannah Rashkin, Eric Michael Smith, Margaret Li, and Y-Lan Boureau. 2019. Towards empathetic open-domain conversation models: A new benchmark and dataset. In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 5370–5381, Florence, Italy. Association for Computational Linguistics.	960
918		961
919		962
920		963
921		964
922		965
923		966
924	Nils Reimers and Iryna Gurevych. 2019. SentenceBERT: Sentence embeddings using Siamese BERT-networks. In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing</i> , pages 3982–3992, Hong Kong, China. Association for Computational Linguistics.	967
925		968
926		969
927		970
928		971
929		972
930		973
931		974
932	Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2019. WinoGrande: An adversarial Winograd Schema Challenge at scale. In <i>Proceedings of the Thirty-Fourth AAAI Conference on Artificial Intelligence</i> , New York City, USA. AAAI Press.	975
933		976
		977
		978
		979
		980
		981
		982
		983
		984
		985
		986
		987
		988
		989
		990
		991
		992
		993
		994
		995
		996
		997
		998
		999
		1000

990 with perfect voice assistants. In *Proceedings of the*
991 *2021 CHI Conference on Human Factors in Comput-*
992 *ing Systems*, CHI '21, pages 1–15, New York, NY,
993 USA. Association for Computing Machinery.

994 Traci Walker, Paul Drew, and John Local. 2011.
995 [Responding indirectly](#). *Journal of Pragmatics*,
996 43(9):2434–2451.

997 Jason Williams, Antoine Raux, Deepak Ramachandran,
998 and Alan Black. 2013. [The Dialog State Tracking](#)
999 [Challenge](#). In *Proceedings of the SIGDIAL 2013*
1000 *Conference*, pages 404–413, Metz, France. Associa-
1001 tion for Computational Linguistics.

1002 Rowan Zellers, Yonatan Bisk, Roy Schwartz, and Yejin
1003 Choi. 2018. [SWAG: A large-scale adversarial dataset](#)
1004 [for grounded commonsense inference](#). In *Proceed-*
1005 *ings of the 2018 Conference on Empirical Methods in*
1006 *Natural Language Processing*, pages 93–104, Brus-
1007 sels, Belgium. Association for Computational Lin-
1008 guistics.

1009 Houyu Zhang, Zhenghao Liu, Chenyan Xiong, and
1010 Zhiyuan Liu. 2020. [Grounded Conversation Gener-](#)
1011 [ation as Guided Traverses in Commonsense Knowl-](#)
1012 [edge Graphs](#). In *Proceedings of the 58th An-*
1013 *ual Meeting of the Association for Computational*
1014 *Linguistics*, pages 2031–2043, Online. ArXiv:
1015 1911.02707.

1016 Saizheng Zhang, Emily Dinan, Jack Urbanek, Arthur
1017 Szlam, Douwe Kiela, and Jason Weston. 2018. [Per-](#)
1018 [sonalizing dialogue agents: I have a dog, do you have](#)
1019 [pets too?](#) In *Proceedings of the 56th Annual Meet-*
1020 *ing of the Association for Computational Linguistics*,
1021 pages 2204–2213, Melbourne, Australia. Association
1022 for Computational Linguistics.

1023 Dengyong Zhou, Qiang Liu, John Platt, and Christopher
1024 Meek. 2014. Aggregating ordinal labels from crowds
1025 by minimax conditional entropy. In *Proceedings of*
1026 *the 31st International Conference on Machine Learn-*
1027 *ing*, pages 262–270, Beijing, China. ACM Press.

1028 Pei Zhou, Karthik Gopalakrishnan, Behnam Hedayat-
1029 nia, Seokhwan Kim, Jay Pujara, Xiang Ren, Yang
1030 Liu, and Dilek Hakkani-Tur. 2022. [Think before](#)
1031 [you speak: Explicitly generating implicit common-](#)
1032 [sense knowledge for response generation](#). In *Pro-*
1033 *ceedings of the 60th Annual Meeting of the Asso-*
1034 *ciation for Computational Linguistics*, pages 1237–
1035 1252, Dublin, Ireland. Association for Computational
1036 Linguistics.

A Ethical Considerations

Undesired bias and abusive content: A multitude of sources have reported that data-driven conversational systems can (re)produce undesired bias or abusive language existing in language resources used for development. To minimize such a risk, we carefully curated conversation examples in SUGAR. Our target task is response selection, where systems only produce language in a pre-compiled response list, and therefore, it is not likely that resulting systems yield harmful content. However, users of SUGAR should be cautious when it is used for developing generation systems in future work.

Human subjects: Crowd workers in Amazon Mechanical Turk (MTurk) participated in our data collection pipeline. Our annotation tasks were reviewed by the institutional review process before being published in MTurk to avoid ethical issues. We did not collect any personally identifiable information of workers other than (anonymized) Turker IDs. Task rewards were decided by several rounds of trials so that workers can receive at least \$6.50 hourly.

Use of external data and tools: We used external datasets such as ATOMIC₂₀ and ConceptNet and tools such as spaCy and Transformers library. We have confirmed that the use of these resources for our research does not violate usage restrictions.

B Manual Annotation

We recruited non-expert crowd workers in Mturk in annotation steps (2-5). In all steps, crowd workers were required to meet the following qualification requirements: (i) Their number of tasks approved $\geq 5k$, (ii) the task approval rate $\geq 99\%$, (iii) their location is the US, and (iv) they answer an exercise question correctly. Figure 5 shows our annotation interfaces for crowdsourcing.

Two of the authors were involved in the annotation steps (1), (4), (5), and (8). They are ESL with a degree in computer science from a school in the US (one holds a master’s degree, and the other holds a Ph.D.). They all have backgrounds in NLP/CL research.

C Distractor Selection

This section presents the technical details of the distractor selection method (Step 7). Below, tunable

parameters like thresholds on scores and the number of iterations were empirically selected based on several pilot runs.

C.1 Response Selection

Our method selects distractor responses from all the responses in the dataset in two steps: We first create an initial dataset by a light-weight method (Algorithm 1) and then perform adversarial filtering (Algorithm 2).

First Step (Algorithm 1)

The objective of the first step is to avoid including false-negative responses (Lines 3-6). We discard responses that are too similar to r_1 in terms of the overlap coefficient of content words (noun, verb, adjective, and adverb).

$$\text{Overlap}(x, y) = \frac{|\text{CW}(x) \cap \text{CW}(y)|}{\min(|\text{CW}(x)|, |\text{CW}(y)|)}, \quad (1)$$

where $\text{CW}(x)$ is a set of content words in x . We set the threshold of overlap coefficient to 0.75. We use the same constraint on their goal texts. We also measure their closeness by the cosine similarity of their sentence embeddings (denoted as EmbSim) and discard candidates whose similarity is 0.5 or higher. We then sample $m - 1$ responses from this filtered response pool one by one (Lines 11-15). To diversify response options, we remove similar responses to the picked one from the pool based on overlap coefficient (Line 16-19).

Second Step (Algorithm 2)

We then perform $J = 3$ rounds of adversarial filtering. Our method is a slightly modified version of the algorithm used by Bhagavatula et al. (2020). In each round, we split the dataset into $K = 10$ folds (Line 6), and for each split, we train a binary logistic regression classifier that takes sentence embeddings of u , S_1 , and a response candidate $r \in R$ (Line 8). We pre-compute their sentence embeddings with the pre-trained Sentence-Transformers (Reimers and Gurevych, 2019) with MPNet (Song et al., 2020). Once the classifier is trained, we score response candidates in each example and identify distractors whose scores are lower than that of the reference response r_1 plus a margin $= 0.05$. We replace these *easy* distractors with more confusing ones (Line 14-16). In this way, we repeatedly update the dataset (Line 17) and output the final result (Line 18).

Algorithm 1 Create an initial dataset by light-weight filtering

Input: m , Dataset $\mathcal{D} = \{(u^{(i)}, g^{(i)}, r_1^{(i)}, S_1^{(i)})\}_{i=1, \dots, N}$, $\triangleright N :=$ number of examples in the dataset.
Output: $\mathcal{D}' = \{(u^{(i)}, g^{(i)}, R^{(i)}, S_1^{(i)})\}_{i=1, \dots, N}$ $\triangleright R^{(i)} := \{r_1^{(i)}, \dots, r_m^{(i)}\}$ \triangleright Initial dataset

```
1: function INITDATASET( $m, \mathcal{D}$ )
2:    $\mathcal{D}' \leftarrow \emptyset$ 
3:   for  $i : 1..N$  do
4:      $\mathcal{P} \leftarrow \{r_1^{(j)}\}_{j=i, \dots, i-1, i+1, \dots, N}$   $\triangleright$  All the responses in  $\mathcal{D}$  but  $r_1^{(i)}$ 
5:     # (1) Remove too similar responses
6:     for  $j : 1..N$  do
7:       if  $i=j$  then
8:         continue
9:       if  $\text{Overlap}(u^{(i)}, u^{(j)}) \geq 0.75$  or  $\text{Overlap}(g^{(i)}, g^{(j)}) \geq 0.75$ 
10:        or  $\text{EmbSim}(u_1^{(i)}, r_1^{(j)}) \geq 0.5$  then
11:          Remove  $r_1^{(j)}$  from  $\mathcal{P}$ 
12:        # (2) Pick  $m-1$  similar responses
13:         $R^{(i)} \leftarrow \{r_1^{(i)}\}$ 
14:        for  $j : 1..m-1$  do
15:          Sample  $r \in \mathcal{P}$ 
16:          Add  $r$  to  $R^{(i)}$ 
17:        # (3) Remove similar responses from the pool
18:        for all  $r' \in \mathcal{P}$  do
19:          if  $\text{Overlap}(r, r') \geq 0.75$  then
20:            Remove  $r'$  from  $\mathcal{P}$ 
21:        Add  $(u^{(i)}, g^{(i)}, R^{(i)}, S_1^{(i)})$  to  $\mathcal{D}'$ 
22:   return  $\mathcal{D}'$ 
```

C.2 Situation Selection

Next, we update S_1 , which only contains relevant information to u and r_1 , to include l statements in total such that some of them are associated with distractors or not directly related to the conversation. Otherwise, reference responses can be easily identified by superficial clues. Having irrelevant situation statements is also for simulating real use cases, where a conversational system has access to a wide range of sensory information or external APIs, but most of them are unimportant for addressing a user’s request.

It is required that (a) additional situation statements do not disqualify the reference response, and (b) they do not contradict others. To this end, we again use sentence embeddings with keyword-based heuristics. We first combine the statements associated with distractor responses and create a pool of candidates. Here, we drop statements that are similar to the response candidates in terms of the overlap coefficient of content words with a threshold of 0.75. We also used manually defined keywords to discard situation statements that tend

to contradict others (e.g., the time is midnight, the user is injured, etc.). We then iterate over six categories and pick situation statements from the pool one by one. We score statement s of category c using the function below:

$$f(s; R, S') = \max_{r \in R} \text{EmbSim}(s, r) \quad (2)$$

$$- \max_{s' \in S'_c} \text{EmbSim}(s, s') \quad (3)$$

$$- \frac{1}{2} \max_{s' \in S'_c \setminus \{c\}} \text{EmbSim}(s, s'), \quad (4)$$

where S' is the current situation statements, $S'_c \subset S'$ represents the statements in S' of category c , and \mathcal{C} denotes a set of situation categories. We pick distractor statements until we exhaust all the candidates in the pool or the maximum score does not reach 0. We then draw statements from the entire dataset in the same way until $|S|$ reaches $l = 12$. For time, date, behavior, and location categories, we pick zero or one statement as those categories are not likely to have more than one value.

Algorithm 2 Adversarial filtering (AF) for R

Input: m , Dataset $\mathcal{D} = \{(u^{(i)}, g^{(i)}, r_1^{(i)}, S_1^{(i)})\}_{i=1, \dots, N}$, $\triangleright N :=$ number of examples in the dataset.
Output: $\mathcal{D}' = \{(u^{(i)}, g^{(i)}, R^{(i)}, S_1^{(i)})\}_{i=1, \dots, N}$ $\triangleright R^{(i)} := \{r_1^{(i)}, \dots, r_m^{(i)}\}$
1: $\mathcal{P} \leftarrow \{(r_0)_i\}$ \triangleright All responses in \mathcal{D}
2: (1) Create an initial dataset D_0
3: $\mathcal{D}_0 \leftarrow \text{INITDATASET}(m, \mathcal{D})$ \triangleright See Algorithm 1
4: (2) Run AF for J rounds
5: **for** $j : 1..J$ **do** \triangleright We set $J = 3$
6: Split \mathcal{D}_{j-1} into K -folds $\{(\mathcal{T}^k, \mathcal{V}^k)\}_{k=1, \dots, K}$ \triangleright We set $K = 10$
7: **for** $k : 1..K$ **do**
8: Train a binary logistic regression classifier \mathcal{M} on \mathcal{T}^k
9: **for all** $(u, g, R, S_1) \in \mathcal{V}^k$ **do**
10: **for all** $r \in R \setminus \{r_1\}$ **do**
11: (f : \mathcal{M} 's score function)
12: **if** $f(r) + \gamma \leq f(r_1)$ **then** $\triangleright \gamma$ is a margin, which we set to 0.05.
13: Remove r from R
14: Pick r' s.t. $f(r') - \gamma > f(r_1)$
15: Add r' to R
16: Update \mathcal{V}^k with the new R
17: $\mathcal{D}_j \leftarrow \bigcup_{k=1}^K \mathcal{V}_k$
18: $\mathcal{D}' \leftarrow \mathcal{D}_K$ \triangleright End

D Situations

We provide the definitions and examples of situation categories in Table 4.

E Examples in SUGAR

Table 5 shows one example in SUGAR.

Exercise

Answer exercise questions below to proceed to a task. Your responses will be rejected if you don't complete the exercise.

User Goal: to go outside

Request: "Can you turn off the lights?" →

Situation: The user is about to go out. It's going to rain today.

Question: Which is more helpful?

Response 1: Sure, I turned off the lights.

Response 2: Sure, please make sure to take an umbrella. It's going to rain today.

Test your answer

(a) In addition to annotation guidelines, we provide one exercise question per task to train crowd workers. We used exercise questions in all the crowdsourced annotation tasks in our pipeline.

Goal	to do exercise
Request	I'd like to drink water.
Response	<p>Sure/Yes,</p> <div style="border: 1px solid #ccc; padding: 5px; margin-bottom: 5px;"> One sentence to make a follow-up request to achieve the goal / suggestion. </div> <p><small>Try to write follow-up requests/additional suggestions to help the user achieve the goal. Examples: "Please turn on AC" -> "Make sure the window is closed", "I'd like some cold water." -> "Would you like ice?" If you cannot come up with any, please write 'none' and skip.</small></p>

(b) In Step 2 (response collection), we instructed workers to write suggestions to achieve the given user goal.

1. The user said:	2. The robot observed that... (Please edit the pre-written text)	3. So,
Can you bring the cake?	<div style="border: 1px solid #ccc; padding: 10px; min-height: 100px;"> [user] is invited to a party. </div> <p>Things you should include:</p> <ol style="list-style-type: none"> 1. When the robot heard the request, the robot guessed that the user wants to go to party because the robot observed ____. 2. The robot accepted or rejected the request because the robot observed ____. 3. The robot gave an additional suggestion/request because the robot observed ____. 	<p>The robot guessed that</p> <div style="border: 1px solid #ccc; padding: 5px; margin-bottom: 5px; text-align: center;"> the user wants to go to party. </div> <p>and said</p> Yes, I got it. Please remember to bring the gift as well.

(c) In Step 3 (situation collection), we guided workers to write observable facts for filling the gap among the provided user utterance, goal, and response.

Figure 5: Annotation interfaces for Step 2 (response collection) and Step 3 (situation collection)

Category	Definition	Example
Location	Information about the user's current location.	The user is home. / The user is at the entrance of a house.
Possession	Information about what the user possesses.	The user owns a car. / There are apples in the kitchen.
Time	Information about time.	It's midnight. / It's morning.
Date	Information about date and season.	It's the user's birthday. / It's summer.
Behavior	Information about the user's behavior.	The user has just woken up. / The user came back from jogging.
Environment	Information about non-user entities (person, objects, etc.).	The room temperature is hot. / The user's car has a flat tire.

Table 4: Definitions of the situation categories.

Utterance	Please turn on the TV.
Situations	<p>It is evening now.</p> <p>[user] is home.</p> <p>[user] is in the living room.</p> <p>[user] is sitting on the couch.</p> <p>[user] has a TV in the house.</p> <p>[user] has an outfit on the bed.</p> <p>[user] has drinks and snacks in the kitchen.</p> <p>[user] has game cards on the shelf.</p> <p>The TV is off.</p> <p>[someone]'s birthday is today.</p> <p>There are several sports games available to watch.</p> <p>There is a basketball game scheduled.</p>
Responses	<p>Sure. Would you like me to check today's sports listings? (<i>Best</i>)</p> <p>Sure. Shall I pour a drink and bring some snacks for the game? (<i>Acceptable</i>)</p> <p>Sure, shall I select an outfit for you? (<i>Bad</i>)</p>

Table 5: Response selection example in SUGAR