

Knowledge-Guided Interpretation and Generation of Task-Oriented Dialogue

Alfredo Gabaldon, Pat Langley, Ben Meadows and Ted Selker

1 **Abstract** In this paper, we present an architecture for task-oriented dialogue that
2 integrates the processes of interpretation and generation. We analyze implemented
3 systems based on this architecture—one for meeting support and another for assisting
4 military medics—and discuss results obtained with the first. In closing, we review
5 some related dialogue architectures and outline plans for future research.

6 1 Introduction

7 Systems that use natural language to assist a user in carrying out some task must
8 interact with that user as execution of the task progresses. The system in turn must
9 interpret the user's utterances and other environmental input to build a model of what
10 both it and the user *believe* and *intend*—in regard to each other and the environment.
11 The system also requires knowledge to use the model it constructs to participate in
12 a dialogue with the user and support him in achieving his goals.

A. Gabaldon (✉) · P. Langley · T. Selker
Silicon Valley Campus, Carnegie Mellon University, Moffett Field, CA 94035, USA
e-mail: alfredo.gabaldon@ge.com

A. Gabaldon
Present Address: GE Global Research, 1 Research Circle, Niskayuna, NY 12309, USA

P. Langley
Present Address: Department of Computer Science, University of Auckland,
Auckland 1142, New Zealand
e-mail: patrick.w.langley@gmail.com

B. Meadows
Department of Computer Science, University of Auckland, Auckland 1142, New Zealand
e-mail: bmea011@auckland.ac.nz

T. Selker
Present Address: Visiting Associate Professor Aarhus University, Aabogade 34, DK-8200
Aarhus, Denmark
e-mail: ted.selker@gmail.com

© Springer International Publishing Switzerland 2016

A. Rudnicky et al. (eds.), *Situated Dialog in Speech-Based*

Human-Computer Interaction, Signals and Communication Technology,

DOI 10.1007/978-3-319-21834-2_3

13 In this paper we report on two systems we have built for task-oriented dialogue
14 and describe the architecture that underpins them. The architecture integrates two
15 processes: *dialogue interpretation*, which builds an expanding model of the user's
16 context in terms of their beliefs and goals, and *dialogue generation*, which uses this
17 interpretation of the situation and background knowledge to help the user achieve
18 his goals through a task-directed conversation.

19 In addition to integrating interpretation and generation, the architecture incor-
20 porates several other important features. Although we believe that domain-specific
21 knowledge is essential in intelligent systems, we also believe that intelligent behavior
22 relies on abstract *meta-level* knowledge that generalizes across different domains. In
23 particular, we are interested in high-level aspects of dialogue: knowledge and strate-
24 gies relevant to dialogue processing that are independent of the actual content of
25 the conversation. The architecture separates domain-level from meta-level content,
26 using both during interpretation and generation. The work we report is informed
27 by cognitive systems research, a key feature of which is arguably integration and
28 processing of knowledge at different levels of abstraction [7].

29 Another feature of our architecture is the incremental nature of its processes. We
30 assume that dialogues occur within a changing environment and that the tasks to
31 be accomplished are not predetermined but discerned as the dialogue proceeds. Our
32 architecture incrementally expands its understanding of the situation and the user's
33 goals, acts according to this understanding, and adapts to changes in the situation,
34 sometimes choosing to pursue different goals. In other words, the architecture sup-
35 ports *situated* systems that carry out *goal-directed* dialogues to aid their users. In
36 the next section we discuss two implemented prototypes that demonstrate this key
37 functionality. We follow this with a detailed description in Sect. 3 of the underlying
38 architecture and a discussion of results in Sect. 4. We conclude with comments on
39 related work and plans for future research.

40 2 Application Systems for Task-Oriented Dialogue

41 In this section, we discuss two prototypes that incorporate our architecture as their
42 dialogue engine. The first system facilitates cyber-physical meetings by interacting
43 with humans and equipment; the second is an advisory system that collaborates with
44 a military medic to address the mutual goal of treating a patient. In each case, we
45 discuss the setting, the knowledge that drives behavior, and a sample dialogue that
46 illustrates operation. Neither system interacts with users in spoken language, but
47 each could be adapted to this end in a straightforward manner.

48 **2.1 A System for Meeting Support**

49 We have developed a system that supports ‘cyber-physical’ meetings in which a
 50 number of human users interact with each other and with a number of devices to
 51 pursue some joint activity. The specific scenario we examine here involves a session
 52 between a physical therapist, a patient, and a relative of the patient who have the
 53 shared aim of completing a physical therapy session. During the meeting, the users
 54 carry out a dialogue that includes the system as a participant. The users and system
 55 also interact with distributed sensors that measure leg motion, heart rate, and ambient
 56 light and temperature, as well as distributed actuators that include a ‘smart’ television
 57 and a text-to-speech generator.

58 During a physical therapy session, users communicate by entering sentences
 59 through a menu-based smart-phone interface, although we intend to replace this with
 60 a spoken-language interface in future work. The current interface accepts sentences
 61 in a controlled version of English, translates them into a logic-like representation,
 62 and sends the result to the dialogue system for processing. At the same time, the
 63 sensing devices send messages to the dialogue system via HTTP. Each user can see
 64 the utterances of other users and those of the system via the phone interface.

65 The system has knowledge about activities that arise during a therapy session,
 66 including physical activities the patient can carry out, such as repeatedly lifting a
 67 leg, and therapist procedures, such as obtaining relevant patient information (from a
 68 database or the patient), measuring range of movement, and tracking heart rate and
 69 blood pressure. In a typical session, the system asks the patient standard questions
 70 and the therapist asks additional ones. Then the therapist decides on the task for

Table 1 Sample dialogue involving a physical therapist (Doctor), a patient (John), and the meeting support system

Doctor	John, please do leg lifts.
Patient	Okay. I will do leg lifts.
TV	[starts playing tutorial video about leg lifting]
Sensor	[sends movement signal to system]
System	Great start John!
Sensor	[sends movement signal to system] [after patient does a few more leg lifts ...]
System	You are doing great, John! [after a few more, the patient stops ...]
Patient	I can't do any more leg lifts.
Doctor	It's okay, John. That was good.
System	John, you did eight repetitions.
Patient	How many did I do in the previous session?
System	In the previous session you did five repetitions.

71 the session and suggests it to the patient, relative, and system. If the patient or
72 family member accepts the proposed task, the system updates its model of the current
73 situation and proceeds accordingly. The system supports the patient's execution of a
74 task by tracking his progress, instructing the television to show a tutorial video, and
75 providing feedback. For instance, once sensor input reveals the patient has started
76 doing an exercise, it might encourage him by saying "Great start!"

77 Specific components of the meeting support system include a menu-based inter-
78 face on a smart phone to input English sentences, a phone application that serves
79 as a motion detector, a television for displaying tutorials and other support videos,
80 a heart-rate monitor, environmental sensors for temperature and lighting, an HTTP
81 client/server module for component communication, and the dialogue system. Table 1
82 shows a sample dialogue for one of the physical therapy scenarios. In this case, the
83 patient John participates in a session in which he partially complete a leg exercise
84 under supervision of a therapist at a remote location. We will return to this case study
85 in Sect. 4, where we examine it in more detail.

86 **2.2 A Medic Assistant**

87 Our second prototype involves scenarios in which a military medic on the battlefield
88 helps an injured teammate. Because the medic has limited training, he interacts with
89 the dialogue system to get advice on treating the person; the system plays the role of a
90 mentor with medical expertise. The medic and system collaborate towards achieving
91 the shared goal of stabilizing the patient's medical condition. The system does not
92 know the specific task in advance. Only after the conversation starts, and the medic
93 provides relevant information, does the system act on this content and respond in
94 ways that are appropriate to achieving the goal. The system does not effect change
95 on the environment directly; the medic provides both sensors and effectors, with the
96 system influencing him by giving instructions.

97 During an interaction, the system asks an initial sequence of questions that lead
98 the medic to provide details about the nature of the injury. This sequence is not pre-
99 determined, in that later questions are influenced by the medic's responses to earlier
100 ones. Table 2 shows a sample dialogue in which the medic-system team attempts
101 to stabilize a person with a bleeding injury. The system possesses domain knowl-
102 edge about how to treat different types of injuries, taking into account their location,
103 severity, and other characteristics. The program can also adapt the treatment accord-
104 ing to the medic's situation. For instance, it may try a different treatment for a wound
105 if the medic claims that he cannot apply a particular treatment because he lacks the
106 supplies necessary for that purpose.

107 This system uses a Web interface similar to a text-messaging application, although
108 again we plan to replace this with a spoken dialogue module in the future. The medic
109 types English sentences into a form element within the interface, which it sends to
110 the dialogue system via an HTTP request. The system in turn sends the content to a
111 natural language processor that translates into a logical form our system can interpret.

112 We have used *Skyphrase* (<http://www.skyphrase.com>), a proprietary, semantically-
 113 driven language interpreter, for this purpose. Skyphrase sends back the translation,
 114 again via HTTP, and the dialogue system uses the information to update its model
 115 of the current situation. This drives the system’s continuing behavior in pursuing the
 116 joint task. Lack of space constrains our discussion of the system here, but we have
 117 provided further details elsewhere [6].

118 **2.3 Discussion**

119 Although the two systems we have just described are limited in many respects,
 120 developing them raised a number of new challenges that have not been fully addressed
 121 by existing dialogue systems, at least in combination:

- 122 • As the dialogue and the users’ situation evolves, the system must update its account
 123 of the situation by interpreting new information as it arrives, including dialogue
 124 utterances and sensor inputs.
- 125 • The system’s interaction with users is goal directed and involves joint activity
 126 over time; this activity includes not only actions carried out by the users, but also
 127 communicative actions and commands to device actuators.
- 128 • The system must establish a *common ground* [4] with its users, that is, it must
 129 develop a shared model of the environmental situation and others’ mental states; it
 130 must also update this joint model as circumstances change and information flows
 131 among participants.

Table 2 Sample dialogue between a human medic and a more knowledgeable computer advisor

Medic	We have a man injured!
Advisor	Okay. What type of injury?
Medic	He’s bleeding.
Advisor	How bad is the bleeding?
Medic	Pretty bad. I think it is an artery.
Advisor	Okay. Where is the injury?
Medic	It’s on the left leg.
Advisor	Apply pressure on the leg’s pressure point.
Medic	Roger that.
Advisor	Has the bleeding stopped?
Medic	No. He’s still bleeding.
Advisor	Okay. Apply a tourniquet.
Medic	Where do i put the tourniquet?
Advisor	Just below the joint above the wound.
Medic	Okay. The bleeding has stopped.

- 132 • The beliefs and goals of each participant are not stated explicitly, but the system
133 must infer enough of them to be effective; this involves using not only domain-
134 specific knowledge, but also more abstract knowledge that relates mental states to
135 communication events.
- 136 • The overall process is highly dynamic, as the system continuously draws inferences
137 from users' utterances and other input to expand its understanding of the evolving
138 situation, and as it carries out activities to achieve goals as they arise.

139 Our application systems and architecture represent first steps towards addressing
140 these challenges. In the next section we describe the integrated architecture, an imple-
141 mentation of which serves as the main component of the two systems above.

142 3 Agent Architecture

143 Now we can turn to our framework for task-oriented dialogue. We have focused on
144 supporting goal-directed behavior that is physically situated in dynamic contexts. The
145 architecture depends on a knowledge base that lets it generate inferences, introduce
146 goals, and execute actions. Input is *multi-modal* in that it might come from speech,
147 text, visual cues, or external sensors. We have implemented the architecture in Prolog,
148 making use of its support for embedded structures and pattern matching, but its
149 representation and control mechanisms diverge substantially from the default Prolog
150 inference engine, as we will see shortly.

151 3.1 Representation and Content

152 As in research on cognitive architectures [9], we distinguish between a dynamic
153 short-term or *working* memory, which stores external inputs and inferences based
154 upon this information, and a more stable *long-term* memory, which serves as a store
155 of knowledge that is used to make inferences and organize activities.

156 Working memory is a rapidly changing set of ground literals that contains the sys-
157 tem's beliefs and goals as it models the evolving situation. Literals for domain-level
158 content, which do not appear as top-level elements in working memory, are stored
159 as relational triples, as in $[il, type, injury]$ or $[il, severity, major]$. This reification
160 lets the system examine and refer separately to different aspects of a single complex
161 concept, including its predicate.

162 Our representation also incorporates meta-level predicates, divorced entirely from
163 the domain level, to denote speech acts [1, 13]. The literature contains many alter-
164 native taxonomies for speech acts; we have adopted a reduced set of six types that
165 has been sufficient for our current purposes. These include:

166 All domain-level and meta-level concepts in working memory are embedded within
167 one of two predicates that denote aspects of mental states: $belief(A, C)$ or $goal(A,$

inform(S, L, C): speaker S asks L to believe content C ;
acknowledge(S, L, C): S tells L it has received and now believes content C ;
question(S, L, C): S asks L a question C ;
propose(S, L, C): S asks L to adopt goal C ;
accept(S, L, C): S tells L it has adopted goal C ;
reject(S, L, C): S tells L it has rejected goal C .

168 C) for some agent A and content C , as in *belief*(*medic*, [*i*, *type*, *injury*]). A mental
 169 state's content may be a triple, [*i*, *r*, *x*], a belief or goal term (*nested* mental states),
 170 an agent's belief that some attribute has a value, as in *belief_wh*(A , [*i*, *r*]), a belief
 171 about whether some propositional content is true, as in *belief_if*(A , C), or a meta-level
 172 literal, such as the description of a speech act.

173 Long-term memory contains generic knowledge in the form of rules. Each rule
 174 encodes a situation or activity by associating a set of triples in its head with a pattern of
 175 concepts in its body. High-level predicates are defined by decomposition into other
 176 structures, imposing an organization similar to that in hierarchical task networks
 177 [11]. Structures in long-term memory include conceptual knowledge, skills, and
 178 goal-generating rules.

179 *Conceptual knowledge* comprises a set of rules which describe classes of situations
 180 that can arise relative to a single agent's beliefs or goals. These typically occur at the
 181 domain level and involve relations among states of the world. Conceptual rules define
 182 complex categories in terms of simpler ones and organize these relational predicates
 183 into taxonomies.

184 *Skills* encode the activities that agents can execute to achieve their goals. Each skill
 185 describes the effects of some action or high-level activity under specified conditions.
 186 The body of a skill include a set of preconditions, a set of effects, and a set of
 187 invariants, along with a sequence of subtasks that are either executable actions, in
 188 the case of primitive skills, or other skills, in the case of nonprimitive skills.

189 *Goal-generating rules* specify domain-level knowledge about the circumstances
 190 under which an agent should establish new goals. For example, an agent might have
 191 a rule stating that, when a teammate is injured, it should adopt a goal for him to be
 192 stabilized. These are similar to conceptual rules, but they support the generation of
 193 goals rather than inference of beliefs.

194 The architecture also includes more abstract, domain-independent knowledge
 195 at the meta-level. This typically involves skills, but it can also specify conceptual
 196 relations (e.g., about transitivity). The most important structures of this type are
 197 *speech act rules* that explain dialogue actions to patterns of agents' beliefs and goals
 198 *without* making reference to domain-level concepts. However, the *content* of a speech
 199 act is instantiated as in any other concept. For example, the rule for an *inform* act is:

$$\begin{aligned}
 \text{inform}(S, L, C) \leftarrow & \text{belief}(S, C), \\
 & \text{goal}(S, \text{belief}(L, C)), \\
 & \text{belief}(S, \text{belief}(L, C)).
 \end{aligned}$$

201 Here S refers to the speaker, L to the listener, and C to the content of the speech act.
 202 Rules for other speech acts take a similar abstract form.

203 Finally, the architecture assumes additional meta-level knowledge in the form
 204 of a *dialogue grammar* that recursively specifies valid patterns of speech acts. For
 205 example, we can decompose a dialogue into a pattern consisting of a speaker S
 206 proposing P to a listener L , followed by L 's acceptance A to S , followed by a dialogue.
 207 To ensure a coherent account of the conversation, the framework includes meta-level
 208 rules that indicate 'conceptual agreement' between the arguments of speech acts;
 209 these ensure that answers to questions are consistent with the agent's beliefs.

210 **3.2 Architectural Processing**

211 Our dialogue architecture uses these structures to operate in dynamic settings, both
 212 interpreting and responding to inputs in terms of its available knowledge and current
 213 model of the situation. Like a traditional cognitive architecture, it operates in cycles
 214 that access relevant knowledge and use it to guide processing. This includes incre-
 215 mentally extending its view of the common ground and its relation to active goals,
 216 then applying skills that are appropriate to achieving those goals. On each cycle, the
 217 architecture invokes a module for dialogue interpretation followed by another for
 218 dialogue generation. We discuss the operation of each of these in turn.

219 **Dialogue Interpretation** The most basic task confronting a dialogue system is to
 220 understand its common ground with other agents. In natural settings many utterances
 221 are elided and others may be misheard, yet it must still construct models of partici-
 222 pants' mental states, making reasonable assumptions about necessary elements that
 223 are missing from working memory. To this end, the architecture's interpretation stage
 224 incorporates a form of abductive inference. This abduction mechanism prefers expla-
 225 nations that introduce few assumptions as possible while accounting for many of the
 226 'observations' that arrive through speech acts.

227 The process first attempts to build support for a top-level rule, such as the existence
 228 of a dialogue in the pattern of speech acts, without making any assumptions. If it
 229 cannot derive the rule's head in this manner, then increases the tolerance to one
 230 default assumption, then two, and so forth, continuing until reaching a maximum.
 231 If the interpretation module finds a proof within this limit, then it adds the assumed
 232 elements to working memory, where they become available for use on later rounds.

233 The abduction mechanism incorporates new utterances and other observations
 234 into working memory at the start of each cognitive cycle. Their arrival can lead it
 235 to introduce beliefs and goals for the participating agents as default assumptions,
 236 with dialogue grammar rules building upon speech acts and other conceptual rules
 237 lower in the proof tree. The module can also introduce omitted speech acts, such as
 238 implicit acknowledgements, as default assumptions, which serve as terminal nodes
 239 in the extended explanation.

240 **Dialogue Generation** The architecture must also produce some response to con-
 241 tinue the dialogue, which is the responsibility of a second module. On each cycle, the
 242 first stage in this process inspects the goal-generating rules, finding which ones have
 243 conditions that match against the current contents of working memory, instantiating

244 their arguments, and adding new top-level goals as a result.¹ Next, an execution
 245 stage selects a top-level goal to pursue and finds a skill clause with this goal in its
 246 head and with conditions that match working memory. The module repeats this step
 247 recursively, finding a path down through the skill hierarchy that, if executed, should
 248 help in achieving the top-level goal. Upon reaching a primitive skill, the architecture
 249 instantiates its variables and carries out its associated actions.

250 On the next cycle, the module might select the same top-level goal and repeat
 251 this process, but, typically, the conditions of some skills along the previous path
 252 will no longer be satisfied, so the architecture follows a slightly different route. This
 253 leads the agent to carry out subskills in sequence, much as in the ICARUS cognitive
 254 architecture [8]. The execution process is reactive in that it responds to changes in
 255 the situation, but the influence of top-level goals also provides continuity over time.
 256 The result is hierarchical behavior in which the agent traverses the branches of an
 257 AND tree, in which each terminal node is an executed primitive skill, across multiple
 258 cognitive cycles.

259 The response of the dialogue generation mechanism also varies based on the type
 260 of goal. Abduced goals typically result in the execution of a meta-level skill, say
 261 one to communicate an instruction. On the other hand, goals inferred from goal-
 262 generating rules typically result in the execution of domain-specific skills. Inter-
 263 estingly, meta-level and domain-specific knowledge always interact at some point
 264 during processing. For instance, a domain-specific skill may have a meta-level skill as
 265 one of its subskills, while a generic skill, like one for communicating an instruction,
 266 is eventually instantiated with some domain-specific content.

267 4 Empirical Evaluation

268 As mentioned earlier, we have used the architecture to implement two dialogue sys-
 269 tems, one for meeting support and another for advising medics. We only have space
 270 here to report results of test runs with the first of them. We will use the interaction
 271 in Table 1 to illustrate the structures and processes that arise during the system's
 272 operation. For instance, after the doctor's utterance "John, do leg lifts," the abductive
 273 interpretation module produces a working memory² that contains:

274 *belief(dr, propose(dr, john, [[e1, exercise_type, leg_lift], [e1, agent, john]]))*
belief(john, propose(dr, john, [[e1, exercise_type, leg_lift], [e1, agent, john]]))
goal(dr, [[e1, exercise_type, leg_lift], [e1, agent, john]])
goal(dr, goal(john, [[e1, exercise_type, leg_lift], [e1, agent, john]]))
belief(john, goal(dr, [[e1, exercise_type, leg_lift], [e1, agent, john]]))
belief(john, goal(dr, goal(john, [[e1, exercise_type, leg_lift], [e1, agent, john]])))

¹The abductive inference mechanism can also introduce new top-level goals as default assumptions during its processing.

²For readability, we omit the top level predicate *belief(sys, Content)* and only show the *Content*.

275 In other words, after the utterance, the system believes that both the doctor and
 276 John believe a speech act occurred in which the speaker (doctor) proposes that the
 277 listener (John) does a leg-lifting exercise, that the doctor has the goal that John do
 278 leg lifts, that the doctor has the goal that John adopt the goal of leg lifting, and
 279 that John also believes the doctor has these two goals. Upon entering the dialogue
 280 generation module, the system does not find any goal-generating rules or any skills
 281 with conditions that match. For this reason, it does not produce any new goals or
 282 generate any utterances before it completes the cognitive cycle.

283 The next utterance is John's response "Okay. I will do leg lifts," which indicates
 284 that he accepts the doctor's proposal. The system starts a new cycle, with the first
 285 step using abductive inference to expand its model of the common ground by adding
 286 to working memory:

```

belief(john, accept(john, dr, [[e1, exercise_type, leg_lift], [e1, agent, john]]))
belief(dr, accept(john, dr, [[e1, exercise_type, leg_lift], [e1, agent, john]]))
goal(john, [[e1, exercise_type, leg_lift], [e1, agent, john]])
287 goal(john, belief(dr, goal(john, [[e1, exercise_type, leg_lift], [e1, agent, john]])))
belief(dr, goal(john, [[e1, exercise_type, leg_lift], [e1, agent, john]]))
goal(sys, [[e1, exercise_type, leg_lift], [e1, agent, john]])

```

288 At this point, the system believes that both the doctor and John believe an accept
 289 speech act occurred, that John has adopted the goal of leg lifting and wants the doctor
 290 to believe that he now has this goal, that the doctor believes that John has adopted
 291 the goal, and, since it aims to support the joint task and both parties have adopted
 292 the goal, the system adopts the goal for itself.

In this case, the dialogue generation module matches a goal-generating rule against
 these elements, producing a new goal to command the television to play a physical
 therapy tutorial for the patient:

```

goal(sys, sys_message(tv, leg_tutorial, nil))

```

293 The system acts on this goal during the execution stage, invoking a skill that sends
 294 a command to the television to play the corresponding video.

During the next cognitive cycle, the system does not receive any utterance from
 the human users, but a signal does arrive from the motion detector indicating that the
 patient has lifted his leg. The interpretation module adds this information to working
 memory as the fact:

```

observation(motion, [[ep2, type, leg_lift], [ep2, agent, john]])

```

295 where *ep2* is a new constant that denotes an event of type *leg_lift* whose *agent* is
 296 John. In response, the abductive inference process extends the current explanation
 297 by adding the elements:

298 *belief(sys, [ep2, type, leg_lift])*
belief(sys, [ep2, agent, john])
belief(sys, [e1, current_state, active])
belief(sys, [e1, reps_done, 1])
belief(sys, [e1, last_rep_time, 1382124783.0])

The system now believes that a leg-lifting event is ongoing and that the first lift has occurred, so it adds a time stamp for the last repetition of the activity, as the system's knows that a leg-lifting exercise involves ten repetitions. Goal generation then produces an intention for the system to utter an encouragement to the patient:

goal(sys, support(sys, john, activity_start))

299 The execution process focuses on this goal and carries out a skill that produces the
 300 utterance "Great start John!", which it sends to the text interface, making it available
 301 to everyone involved in the meeting.

302 We lack the space to completely analyze the remaining interaction, but it is impor-
 303 tant to note how the system reacts to divergences from the above sequence of events.
 304 Consider the case in which the doctor instead proposes "John lie down" and in which
 305 John counters "No, I will do leg lifts." As there is no agreed upon goal, in this case
 306 the system does not play the tutorial and instead reminds John of the doctor's goal
 307 by uttering "John, the doctor wants you to lie down."

308 Alternatively, consider a variation in which the interaction starts with the original
 309 utterances by the doctor and John about leg lifts, followed by the tutorial, but in which
 310 no signal arrives from the motion detector. In this case, after some time has passed
 311 without the expected motion signal, the system generates a goal to utter "John, you
 312 should strap on the motion detector" and executes a skill that communicates this
 313 content to the patient.

314 The different interactions illustrate the system's ability to respond appropriately
 315 based on its beliefs about the mental state of the users (e.g., whether they adopted the
 316 same goal) and the environmental situation (e.g., that the patient forgot to wear the
 317 motion detector). The dialogue framework supports such reactive responses within
 318 the broader context of the high-level goals it has adopted.

319 5 Discussion

320 Our architecture and the two systems that utilize it take steps towards robust, task-
 321 oriented dialogue systems, but there are some issues that we have not addressed
 322 fully. We remarked earlier that we plan to replace the text-based interfaces with
 323 spoken language interfaces. That move will come with the additional complication
 324 of uncertainty in the meaning of utterances, but we believe our abductive approach to
 325 incremental inference is well situated to handle this issue. We must increase the scope
 326 of explanations to include hypotheses about the meaning of each utterance, possibly

327 using some measure of uncertainty. We should also introduce the ability to revise
328 faulty assumptions that arise in dialogue misunderstandings, to which abduction also
329 lends itself [10]. At the same time, one motivation for developing the architecture
330 was to support robust cognitive systems. This suggests additional research goals,
331 including the ability to execute skills in parallel and to handle unfamiliar tasks through
332 problem solving. We believe that our architecture's representations and mechanisms
333 could be adapted to other contexts beyond task-oriented dialogue that involve social
334 cognition. Examples include settings in which agents provide help without verbal
335 communication and in which self-interested agents take advantage of ignorance and
336 deception [3].

337 The literature reports a number of advanced dialogue managers. RavenClaw [2]
338 separates from the domain level some domain-independent aspects of dialogue man-
339 agement, including turn taking, timing, and error handling. In contrast, we have
340 focused on domain-independent principles at the abstract level of dialogue knowl-
341 edge. Moreover, RavenClaw emphasizes generation, while our architecture balances
342 interpretation and generation. Our architecture is similar to Collagen [12] in that both
343 utilize hierarchical plan structures and construct models of agents' beliefs during
344 interpretation and generation, but a key difference is that Collagen does not separate
345 meta-level from domain knowledge. Also, despite sharing some high-level assump-
346 tions, our abduction mechanism makes the two frameworks operate quite differently.
347 We should also mention TRIPS [5], an integrated system that carries out dialogues to
348 help users generate plans, drawing on knowledge to interpret user input and generate
349 responses. However, TRIPS was designed for the task of plan creation, while our
350 architecture can support any collaborative task given suitable domain knowledge.

351 6 Concluding Remarks

352 In this paper, we presented an architecture for task-oriented dialogue that integrates
353 interpretation and generation, along with two implemented systems that build on it.
354 We discussed results obtained from runs with the meeting support system, demon-
355 strating how it interprets the current situation and, by combining meta-level and
356 domain-level knowledge, supports users by participating actively in the dialogue
357 and issuing commands to actuators.

358 In addition to integrating processes for dialogue interpretation and generation, the
359 framework provides a clear separation of meta-level content from domain expertise,
360 which we maintain is a desirable feature in a cognitive architecture. These suggest
361 that it can serve as a solid foundation for future research on both dialogue systems
362 and other software agents that interact with humans.

363 **Acknowledgments** This research was supported by Grant N00014-09-1-1029 from the Office
364 of Naval Research and a gift from Ericsson. We thank Chitta Baral, Paul Bello, Will Bridewell,
365 Herb Clark, Tolga K onik, Nimish Radia, Ted Selker, David Stracuzzi, Chihiro Suga, and Richard
366 Weyrauch for discussions that influenced the approach reported here.

367 **References**

- 368 1. Austin JL (1962) How to do things with words. Harvard University Press, Cambridge
 369 2. Bohus D, Rudnicky A (2009) The RavenClaw dialog management framework: architecture
 370 and systems. *Comput Speech Lang* 23(3):332–361
 371 3. Bridewell W, Isaac A (2011) Recognizing deception: a model of dynamic belief attribution.
 372 In: *Advances in Cognitive Systems: papers from the 2011 AAAI fall symposium*, pp 50–57
 373 4. Clark HH (1996) Using language. Cambridge University Press, Cambridge
 374 5. Ferguson G, Allen JF (1998) TRIPS: an integrated intelligent problem-solving assistant. In:
 375 *Proceedings of the 15th National Conference on Artificial Intelligence*, pp 567–572
 376 6. Gabaldon A, Langley P, Meadows B (2013) Integrating meta-level and domain-level knowledge
 377 for interpretation and generation of task-oriented dialogue. In: *Proceedings of the Second*
 378 *Annual Conference on Advances in Cognitive Systems*
 379 7. Langley P (2012) The cognitive systems paradigm. *Advances Cognitive Systems* 1:3–13
 380 8. Langley P, Choi D, Rogers S (2009) Acquisition of hierarchical reactive skills in a unified
 381 cognitive architecture. *Cognitive Systems Research* 10:316–332
 382 9. Langley P, Laird JE, Rogers S (2009) Cognitive architectures: research issues and challenges.
 383 *Cognitive Systems Research* 10:141–160
 384 10. McRoy S, Hirst G (1995) The repair of speech act misunderstandings by abductive inference.
 385 *Computational Linguistics* 21(4):435–478
 386 11. Nau DS, Cao Y, Lotem A, Munoz-Avila A (2001) The SHOP planning system. *AI Magazine* 22:
 387 91–94
 388 12. Rich C, Sidner CL, Lesh N (2001) Collagen: Applying collaborative discourse theory to human-
 389 computer interaction. *AI Magazine* 22:15–25
 390 13. Searle J (1969) *Speech acts: An essay in the philosophy of language*. Cambridge University
 391 Press, New York

Author Queries

Chapter 3

Query Refs.	Details Required	Author's response
	No queries.	

UNCORRECTED PROOF