

Robotic assistance in the coordination of patient care

The International Journal of
Robotics Research
1–17

© The Author(s) 2018

Reprints and permissions:

sagepub.co.uk/journalsPermissions.nav

DOI: 10.1177/0278364918778344

journals.sagepub.com/home/ijr



**Matthew Gombolay¹, Xi Jessie Yang¹, Bradley Hayes¹, Nicole Seo¹,
Zixi Liu¹, Samir Wadhwanian¹, Tania Yu¹, Neel Shah², Toni Golen², and Julie Shah¹**

Abstract

We conducted a study to investigate trust in and dependence upon robotic decision support among nurses and doctors on a labor and delivery floor. There is evidence that suggestions provided by embodied agents engender inappropriate degrees of trust and reliance among humans. This concern represents a critical barrier that must be addressed before fielding intelligent hospital service robots that take initiative to coordinate patient care. We conducted our experiment with nurses and physicians, and evaluated the subjects' levels of trust in and dependence upon high- and low-quality recommendations issued by robotic versus computer-based decision support. The decision support, generated through action-driven learning from expert demonstration, produced high-quality recommendations that were accepted by nurses and physicians at a compliance rate of 90%. Rates of Type I and Type II errors were comparable between robotic and computer-based decision support. Furthermore, embodiment appeared to benefit performance, as indicated by a higher degree of appropriate dependence after the quality of recommendations changed over the course of the experiment. These results support the notion that a robotic assistant may be able to safely and effectively assist with patient care. Finally, we conducted a pilot demonstration in which a robot-assisted resource nurses on a labor and delivery floor at a tertiary care center.

Keywords

Human–robot teaming, human–robot interaction, planning and scheduling, situational awareness, workload, preference scheduling

1. Introduction

Service robots are being increasingly utilized across a wide spectrum of clinical settings. They are deployed to improve operational efficiency by delivering and preparing supplies, materials and medications (Bloss, 2011; DiGiose, 2013; Hu et al., 2011; Mutlu and Forlizzi, 2008; Özkil et al., 2009). Existing systems exhibit robust, autonomous capabilities for navigating from point to point while avoiding obstacles (Murai et al., 2012a,b), and initial concerns regarding physical safety around people have largely been addressed. However, these robots are not yet well-integrated into the healthcare delivery process: they do not operate with an understanding of patient status and needs, and must be explicitly tasked and scheduled. This can impose a substantial burden upon the nurse in charge of resource allocation, or “resource nurse,” particularly within fast-paced hospital departments, such as the emergency or labor and delivery units.

Resource nurses are essentially solving an NP-hard (Bertsimas and Weismantel, 2005) problem on the fly: they assign resources such as beds (e.g. for triage, in-patient,

recovery, and operating rooms) while subject to upper- and lower-bound temporal constraints on availability and considering stochasticity in the timing of patient progression from one bed type to another. They must also pair patients with staff nurses, equipment, and resources. The resource nurse’s job is made feasible because staff nurses understand patients’ statuses and needs and will take initiative to accomplish some tasks without being explicitly directed.

As the number and types of hospital service robots increases, these robots must similarly take initiative to provide a net productivity benefit. The need to explicitly task many service robots may degrade the performance of a resource nurse (Chen et al., 2011; Cummings and Guerlain, 2007; Olsen and Wood, 2004), which has implications for

¹Massachusetts Institute of Technology, Cambridge, MA, USA

²Beth Israel Deaconess Medical Center, Boston, MA, USA

Corresponding author:

Matthew Gombolay, Georgia Institute of Technology, North Ave NW, Atlanta, GA 30332, USA

Email: Matthew.Gombolay@cc.gatech.edu

both patient safety and the well-being of healthcare professionals (Brandenburg et al., 2015; Kehle et al., 2011; Pizer and Prentice, 2011; Shipman and Sinsky, 2013).

On the other hand, a robot that autonomously takes initiative when performing tasks may make poor decisions in the absence of oversight. Furthermore, decades of research in human factors cautions against fully autonomous decision making, as it contributes to poor human situational awareness and degradation in the human supervisor's performance (Kaber and Endsley, 1997; Parasuraman et al., 2000; Sheridan, 2011; Wickens et al., 2010). When integrating machines into human cognitive work flows, an intermediate level of autonomy is preferred (Kaber and Endsley, 1997; Wickens et al., 2010), in which the system provides suggestions to be accepted or modified by a human supervisor. Such a system would fall within the "4–6" range on the 10-point scale of Sheridan's levels of automation (Parasuraman et al., 2000).

In prior work (Gombolay et al., 2016b), we investigated the human factors implications of fielding hospital service robots that necessarily reason about which tasks to perform and when to perform them. In particular, we investigated trust in and dependence upon robotic decision support among nurses and doctors on a labor and delivery floor. Studies of human–automation interaction in aviation, another safety-critical domain, have shown that human supervisors can inappropriately trust in and rely upon recommendations made by automation systems (Dismukes et al., 2007). For example, numerous aviation incidents have been attributed to human overreliance on imperfect automation (Dismukes et al., 2007). Other studies have examined the effects of changes in system reliability, and found that this led to suboptimal control allocation strategies and reduced levels of trust in the relevant systems (Desai et al., 2013, 2012). There is also evidence that suggestions provided by embodied agents engender over-trust and inappropriate reliance (Robinette et al., 2016). This concern represents a critical barrier to fielding intelligent hospital service robots that take initiative to participate with nurses in decision making.

This paper presents three novel contributions to the fields of robotics and healthcare. First, through human subject experimentation with physicians and registered nurses, we conducted the first known study involving experts working with an embodied robot on a real-world, complex decision-making task comparing trust in and dependence upon robotic- versus computer-based decision support. Previous studies have focused on novice users and/or simple laboratory decision tasks (Bainbridge et al., 2011; de Visser et al., 2012; Kiesler et al., 2008; Leyzberg et al., 2014). Our findings provide the first evidence that experts performing decision-making tasks appear to be less susceptible to the negative effects of support embodiment, as trust assessments were similar under both the computer-based and robotic decision support conditions. Furthermore, embodiment yielded performance gains compared

with computer-based support after the quality of recommendations changed over the course of the experiment. This provides encouraging evidence that intelligent service robots can be safely integrated into the hospital setting.

Second, decision support generated through action-driven learning from expert demonstration produced high-quality recommendations accepted by nurses and physicians at a compliance rate of 90%. This indicates that a hospital service robot may be able to learn context-specific decision strategies and apply them to make reasonable suggestions for which tasks to perform and when. We note that the learning model was presented in prior work (Gombolay et al., 2016a). This paper provides a novel demonstration of the algorithm, providing evidence that machine learning can be used to effectively learn to emulate the decision-making process in this healthcare domain.

Finally, based on the previous two findings, we conducted the first test demonstration in which a robot-assisted resource nurse on a labor and delivery floor in a tertiary care center. Our robot used machine learning computer vision techniques to read the current status of the labor floor and make suggestions about resource allocation, and used speech recognition to receive feedback from the resource nurse. To the best of the authors' knowledge, this is the first investigation to field a robotic system in a hospital to aid in the coordination of resources required for patient care.

We extend our prior work (Gombolay et al., 2016b) in the following ways. First, we increased our sample size by 41%, from 17 to 24 participants, and report our updated findings. Specifically, we observed stronger evidence that Type I and Type II error rates for the computer-based decision support system are more adversely affected by changes in advice quality than the error rates for a robotic decision support system. Second, we conducted post-hoc analysis exploring the participants' response time (i.e. the time between the end of both the computer-based and robotic decision support systems' spoken recommendations and the response from the participant) as a function of the systems' advice quality and embodiment. Third, we include a fuller presentation of the robotic system fielded at our partner hospital.

2. Background

Whereas the effects of embodiment on engagement in social judgment tasks have been studied extensively and well documented (Bartneck et al., 2009; Kidd and Breazeal, 2004; Kiesler et al., 2008; Takayama and Pantofaru, 2009; Tapus et al., 2009), the relationship between embodiment and humans' levels of trust and dependence is a relatively new area of research (Bainbridge et al., 2011; Kiesler et al., 2008; Leyzberg et al., 2014). This topic is crucial if robots are to become more than companions, but advisors to people.

In this context, trust is defined as “the attitude that an agent will help achieve an individual’s goals in a situation characterized by uncertainty and vulnerability (Lee and See, 2004),” and dependence is a behavioral measure indicating the extent to which users accept the recommendation of robots or virtual agents. Measures of dependence are distinguished according to whether the user makes Type I or Type II errors (Dixon and Wickens, 2006). “Type I” refers to reliance, or the degree to which users accept advice from an artificial agent when it offers low-quality recommendations. “Type II” refers to the extent to which human users reject advice from an artificial agent when the advice is of high quality. The degrees to which a user accepts high-quality advice and rejects low-quality advice are called “appropriate compliance” and “appropriate reliance,” respectively.

Studies examining the effects of embodiment on trust and dependence necessarily include objective assessments of dependence and task performance in addition to subjective assessment of the user’s trust in the system (Bainbridge et al., 2011; Bartneck et al., 2009; de Visser et al., 2012; Kiesler et al., 2008; Leyzberg et al., 2014; Pak et al., 2012).

Bainbridge et al. (2011) conducted an experiment where participants collaborated with either a physical, embodied robot or a video of a robot in a book-moving task (i.e. moving books from one location to another). They found that participants were more willing to comply with unusual requests from the physically-present robot than the robot displayed by video.

Similarly, Leyzberg et al. (2014) conducted experiments investigating the effects of embodiment. In their study, participants were given strategy advice for solving Sudoku-like puzzles. The authors found that embodiment was associated with a higher rate of compliance with advice provided by the robot, and suggested this indicated a greater level of human trust for an embodied robot.

Finally, Kiesler et al. (2008) found that participants consumed fewer calories after receiving health advice from a physically embodied robot, as compared with advice from a video of a robot or an on-screen animated virtual agent.

Studies in human factors and decision support have indicated that increased anthropomorphism also affects user interactions (de Visser et al., 2012; Kulić et al., 2016; Pak et al., 2012). Pak et al. (2012) evaluated how the anthropomorphic characteristics of decision support aids assisting subjects answering questions about diabetes influenced subjective trust and task performance. Their results indicated that younger adults trusted the anthropomorphic decision aid more, whereas older adults were insensitive to the effects of anthropomorphism. Moreover, shorter question response time (after controlling for accuracy) was observed in both age groups, suggesting a performance gain when receiving advice from a more anthropomorphic aid. In another study, de Visser et al. (2012) varied the degree of anthropomorphism of a decision support system while participants performed a pattern-recognition task. The results indicated that the perceived knowledgeable-

ness of the system increased with increasing anthropomorphism; however, their findings on dependence were inconclusive. Kulić et al. (2016) provided a helpful survey of work in developing anthropomorphic agents.

The results from studies with embodied robots must be interpreted with caution because they were primarily focused on situations in which robots produced reliable and high-quality recommendations. There is a growing body of research indicating that the quality of decision support cannot be relied upon, especially during complex tasks (Wickens et al., 2013). Negative consequences of humans blindly depending upon imperfect embodied artificial intelligence have been reported previously: for example, Robi-nette et al. (2016) conducted experiments in which a robot guided human participants during a mock emergency rescue scenario involving a building fire. All participants followed the robot, even when the robot led them down unsafe routes and/or displayed simulated malfunctions and other suspicious behavior.

Although there is ongoing work trying to develop formal methods for safe human–robot interaction (Jansen et al., 2017; Li et al., 2014), imperfect automation persists, and dependence upon such imperfect automation presents serious problems for robotic assistance during safety-critical tasks. This concern is heightened by results from studies indicating increased trust in and reliance upon embodied systems as compared with virtual or computer-based decision support, suggesting a higher possibility of committing Type I errors. However, prior studies on embodiment, trust and dependence were conducted with novices rather than domain experts performing complex real-world tasks. This leaves us with founded concerns but also gaps in our understanding of how human–robot interaction impacts the decision making of expert resource nurses. In the following sections, we describe our experiment and present a positive result for service robots in a hospital setting, with Type I and Type II error rates comparable with those observed for computer-based decision support. Furthermore, embodiment appeared to improve performance, as indicated by a higher degree of appropriate compliance when the quality of advice changed mid-experiment.

3. Experimental investigation

In this section, we describe human-subject experimentation aimed at comparing trust in and dependence upon an embodied robot assistant versus computer-based decision support in a population of physicians and registered nurses. The participants interacted with a high-fidelity simulation of an obstetrics department at a tertiary care center. This simulation provided users the opportunity to assume the roles and responsibilities of a resource nurse, which included assigning labor nurses and scrub technicians to care for patients, as well as moving patients throughout various care facilities within the department.



Fig. 1. An experiment participant pictured receiving advice from the robotic decision support.

We conducted the experiment using a within-subjects design that manipulated two independent variables: *embodiment* (subjects received advice from either a robot or a computer) and *recommendation quality* (subjects received high- or low-quality advice). Each participant experienced four conditions, the quality of advice was blocked and the ordering of the conditions was counterbalanced to mitigate potential learning effects. Figure 1 depicts the experimental setup for the embodied condition.

3.1. Hypotheses and measures

H1 *Rates of appropriate compliance with and reliance upon robotic decision support will be comparable to or greater than those observed for computer-based decision support.* Objective measures of compliance and reliance were assessed based on the participants’ “accept” or “reject” responses to each decision support recommendation. Statistics on appropriate compliance, appropriate reliance, Type I and Type II errors were recorded.

H2 *Robotic decision support will be rated more favorably than computer-based decision support in terms of trust and other attitudinal measures.* Results from numerous studies have demonstrated that embodied and anthropomorphic systems are rated more favorably by users than computer-based interactive systems. We hypothesized that the robotic system in our study would elicit this favorable response (H2), while engendering appropriate rates of compliance and reliance (H1). This would indicate a positive signal for the successful adoption of a hospital service robot that participates in decision making. Subjective measures of trust and attitudinal response were collected via questionnaires administered to each participant after each of the four trials. Trust was assessed by a composite rating of seven-point Likert-scale responses for a commonly used, validated trust questionnaire (Jian et al., 2000). Other attitudinal questions were drawn from work by Lee et al. (2006) to evaluate personality recognition, social responses and social presence in human–robot interaction, and were responded to on a 10-point Likert scale.

3.2. Materials and setup

We conducted our experiments using a high-fidelity simulation of a labor and delivery floor. This simulation had previously been developed through a hospital quality-improvement project as a training tool over a year-long, rigorous design and iteration process that included workshops with nurses, physicians, and medical students to ensure the tool accurately captured the role of a resource nurse. Parameters within the simulation (e.g. patient arrivals, timelines on progression through labor) were drawn from medical textbooks and papers and modified through alpha and beta testing to ensure that the simulation closely mirrored the patient population and nurse experience at our partner hospital.

An Aldebaran Nao robot was employed for the embodied condition (Figure 1). A video of the Nao offering advice to a participant with speech and co-speech gestures is shown in Extension 1. Participants received advice through synthesized speech under both the embodied and computer-based support conditions, using a male voice drawn from the Mary Text-to-Speech System (MaryTTS) (Schröder and Trouvain, 2003). The advice was also displayed as text in an in-simulation pop-up box under both conditions. The subject clicked a button to accept or reject the advice; these buttons were not clickable until the narration was complete, which was independent of whether the advice was provided from the embodied versus the computer-based system.

3.3. Experimental procedure

Twenty-four physicians and registered nurses, recruited via email and word-of-mouth, participated in the experiment (one man and 23 women). This gender ratio is a representative sample of our partner OB/GYN department: $\sim 70\%$ of the attending physician population and 100% of the nurses and resident populations are female. The ratio of female-to-male labor nurses in the United States is approximately 9.5 : 1 (and higher in our state), according to data collected by The Kaiser Family Foundation.¹

First, participants provided consent for the experiment and watched an 8-minute tutorial video describing the labor and delivery floor simulation. The tutorial video is available as Extension 2. Participants were instructed to play the simulation four times, with each iteration lasting 10 minutes, simulating a total of 4 hours on the labor floor. The computer or embodied system would interject during the simulation to offer recommendations on which nurse should care for which patient, as well as on patient room assignments. Participants were asked to accept or reject the advice based on their own judgment. They were not informed whether the robotic or virtual decision support coach was providing high- or low-quality advice. Finally, after each of the four trials, participants were asked to rate their subjective experience via a set of Likert-scale questions.

4. Toward decision support: formulation of the resource nurse's decision-making problem

This section provides a formal representation of the resource nurse's decision-making problem. The following section describes how we implemented the decision support based on this formulation.

A resource nurse must solve a problem of task allocation and schedule optimization with stochasticity in the number and types of patients and the duration of tasks. A task τ_i represents the set of steps required to care for patient i , and each τ_i^j is a given stage of labor for that patient. Stages of labor are related by stochastic lower-bound constraints $W_{\langle \tau_i^j, \tau_x^y \rangle}$, requiring the stages to progress sequentially. There are stochastic time constraints, $D_{\tau_i^j}^{abs}$ and $D_{\langle \tau_i^j, \tau_x^y \rangle}^{rel}$, relating the stages of labor to account for the inability of resource nurses to perfectly control when a patient will move from one stage of labor to the next. Arrivals of τ_i (patients) are drawn from stochastic distributions. The model considers three types of patients: patients scheduled for cesarean section, patients scheduled for labor induction, and unscheduled patients. The sets $W_{\langle \tau_i^j, \tau_x^y \rangle}$, $D_{\tau_i^j}^{abs}$, and $D_{\langle \tau_i, \tau_j \rangle}^{rel}$ are dependent upon patient type.

Labor nurses are modeled as agents with a finite capacity to process tasks in parallel, where each subtask requires a variable amount of this capacity. For example, a labor nurse may generally take care of a maximum of two patients simultaneously. If the nurse is caring for a patient who is "full and pushing" (i.e. the cervix is fully dilated and the patient is actively trying to push out the baby) or in the operating room, the nurse may only care for that patient.

Rooms on the labor floor (e.g. a labor room, an operating room, etc.) are modeled as resources, which process subtasks in series. Agent and resource assignments to subtasks are pre-emptable, meaning that the agent and resource assigned to care for any patient during any step in the care process may be changed over the course of executing that subtask.

In this formulation, ${}^tA_{\tau_i^j}^a \in \{0, 1\}$ is a binary decision variable for assigning agent a to subtask τ_i^j for time epoch $[t, t + 1)$. ${}^tG_{\tau_i^j}^a$ is a continuous decision variable for assigning a certain portion of the effort of agent a to subtask τ_i^j for time epoch $[t, t + 1)$. Here ${}^tR_{\tau_i^j}^r \in \{0, 1\}$ is a binary decision variable for whether subtask τ_i^j is assigned resource r for time epoch $[t, t + 1)$, $H_{\tau_i} \in \{0, 1\}$ is a binary decision variable for whether task τ_i and its corresponding subtasks are to be completed, $U_{\tau_i^j}$ specifies the effort required from any agent to work on τ_i^j , and $s_{\tau_i^j}, f_{\tau_i^j} \in [0, \infty)$ are the start and finish times, respectively, of τ_i^j .

The equations

$$\min fn \left(\{ {}^tA_{\tau_i^j}^a \}, \{ {}^tG_{\tau_i^j}^a \}, \{ {}^tR_{\tau_i^j}^r \}, \{ H_{\tau_i} \}, \{ s_{\tau_i^j}, f_{\tau_i^j} \} \right) \quad (1)$$

$$\sum_{a \in A} {}^tA_{\tau_i^j}^a \geq H_{\tau_i}, \forall \tau_i^j \in \tau, \forall t \quad (2)$$

$$0 = \left({}^tG_{\tau_i^j}^a - U_{\tau_i^j} \right) {}^tA_{\tau_i^j}^a H_{\tau_i}, \forall \tau_i^j \in \tau, \forall t \quad (3)$$

$$\sum_{\tau_i^j \in \tau} {}^tG_{\tau_i^j}^a \leq C_a, \forall a \in A, \forall t \quad (4)$$

$$\sum_{r \in R} {}^tR_{\tau_i^j}^r \geq H_{\tau_i}, \forall \tau_i^j \in \tau, \forall t \quad (5)$$

$$\sum_{\tau_i^j \in \tau} {}^tR_{\tau_i^j}^r \leq 1, \forall r \in R, \forall t \quad (6)$$

$$f_{\tau_i^j} - s_{\tau_i^j} \leq ub_{\tau_i^j}, \forall \tau_i^j \in \tau \quad (7)$$

$$f_{\tau_i^j} - s_{\tau_i^j} \geq lb_{\tau_i^j}, \forall \tau_i^j \in \tau \quad (8)$$

$$s_{\tau_x^y} - f_{\tau_i^j} \geq W_{\langle \tau_i, \tau_j \rangle}, \forall \tau_i, \tau_j \in \tau, \forall W_{\langle \tau_i, \tau_j \rangle} \in \mathbf{TC} \quad (9)$$

$$f_{\tau_x^y} - s_{\tau_i^j} \leq D_{\langle \tau_i, \tau_j \rangle}^{rel}, \forall \tau_i, \tau_j \in \tau, \exists D_{\langle \tau_i, \tau_j \rangle}^{rel} \in \mathbf{TC} \quad (10)$$

$$f_{\tau_i^j} \leq D_{\tau_i^j}^{abs}, \forall \tau_i \in \tau, \exists D_{\tau_i^j}^{abs} \in \mathbf{TC} \quad (11)$$

represent a mixed-integer non-linear program over these variables. We address the objective function in Section 5. A mathematical program solver would seek to minimize the application-specific objective function in Equation (1) subject to the constraints in Equations (2)–(11). Equation (2) enforces that each subtask τ_i^j during each time epoch $[t, t + 1)$ is assigned one agent. Equation (3) ensures that each subtask τ_i^j receives a sufficient portion of the effort of its assigned agent a during time epoch $[t, t + 1)$. Equation (4) ensures that agent a is not oversubscribed.

Equation (5) ensures that each subtask τ_i^j of each task τ_i to be completed (i.e. $H_{\tau_i} = 1$) is assigned one resource r . Equation (6) ensures that each resource r is assigned to only one subtask during each epoch $[t, t + 1)$. Equations (7) and (8) requires the duration of subtask τ_i^j to be less than or equal to $ub_{\tau_i^j}$ and at least $lb_{\tau_i^j}$ units of time, respectively. Equation (9) requires that τ_x^y occurs at least $W_{\langle \tau_i^j, \tau_x^y \rangle}$ units of time after τ_i^j . Equation (10) requires that the duration between the start of τ_i^j and the finish of τ_x^y is less than $D_{\langle \tau_i^j, \tau_x^y \rangle}^{rel}$. Equation (11) requires that τ_i^j finish before $D_{\tau_i^j}^{abs}$ units of time have expired since the start of the schedule.

The stochasticity of the problem arises from the uncertainty in the upper and lower bound of the durations, $(ub_{\tau_i^j}, lb_{\tau_i^j})$, of each of the steps in caring for a patient; the number and types of patients, τ ; and the temporal constraints, \mathbf{TC} , relating the start and finish of each step. These

variables are a function of the resource and staff allocation variables ${}^tR_{\tau_i}^a$ and ${}^tA_{\tau_i}^a$, and patient task state $\Lambda_{\tau_i}^j$, which includes information on patient type (i.e. patients presenting with scheduled induction, scheduled cesarean section, or acute unplanned anomaly), gestational age (i.e. number of days and weeks pregnant), gravida (i.e. number of pregnancies), parity (i.e. number of live births), membrane status (i.e. intact or ruptured), anesthesia status (i.e. whether, and which type of, anesthesia has been administered), cervix status, including dilation (i.e. width of cervical opening), station (i.e. location of baby relative to cervix, depth-wise), and effacement (i.e. degree to which the cervix has decreased in length), time of last exam and, finally, the presence of any co-morbidities. Formally, $(\{ub_{\tau_i}^j, lb_{\tau_i}^j | \tau_i^j \in \tau\}, \tau, \mathbf{TC}) \sim P(\{{}^tR_{\tau_i}^a, {}^tA_{\tau_i}^a, \Lambda_{\tau_i}^j, \forall t \in [0, 1, \dots, T]\})$.

We note that the constraints represented in Equations (2)–(11) are linear except for Equation (3), which is a nonlinear, cubic, equality constraint. As nonlinear constraints are more difficult to solve, it may be desirable to linearize such equations. We can substitute the nonlinear equality constraint in Equation (3) with two linear, non-strict inequality constraints, as shown in

$${}^tG_{\tau_i}^a - U_{\tau_i}^j \leq M \left(2 - {}^tA_{\tau_i}^a - H_{\tau_i} \right), \forall \tau_i^j \in \tau, \forall t \quad (12)$$

$${}^tG_{\tau_i}^a - U_{\tau_i}^j \geq M \left({}^tA_{\tau_i}^a + H_{\tau_i} - 2 \right), \forall \tau_i^j \in \tau, \forall t \quad (13)$$

using the big- M method, where M is a large, positive number. For further information regarding the big- M method, please consult Winston et al. (2003).

4.0.1. Computational complexity. Mathematical programs, such as that outlined in Equations (1)–(11), are typically solved exactly with a branch-and-bound search technique due to the presence of integer variables. The computational complexity of a branch-and-bound search for a solution to the assignment of decision variables to the integer decision variables, ${}^tA_{\tau_i}^a \in \{0, 1\}$, $\{{}^tR_{\tau_i}^r\} \in \{0, 1\}$, and $\{H_{\tau_i}\} \in \{0, 1\}$. Considering a scenario with n tasks, m subtasks per task, a agents (e.g. nurses), r resources (e.g. rooms), and a time horizon of T units of time, the computational complexity for solving the constraint portion of the formulation (i.e. Equations (2)–(11)) would be of the order $O\left(2^{\binom{|{}^tA_{\tau_i}^a| + |{}^tR_{\tau_i}^r| + |H_{\tau_i}|}{\tau_i}}\right) = O(2^{(anmT)(rnmT)(n)}) = O(2^{an^3m^2T^2r})$.

Mathematical programs, such as that outlined in Equations (1)–(11), are typically solved exactly with a branch-and-bound search technique due to the presence of integer variables. The computational complexity of a branch-and-bound search for a solution to the assignment of decision variables would be dominated by searching over

assignments to the integer decision variables, ${}^tA_{\tau_i}^a \in \{0, 1\}$, $\{{}^tR_{\tau_i}^r\} \in \{0, 1\}$, and $\{H_{\tau_i}\} \in \{0, 1\}$. Considering a scenario with n tasks, m subtasks per task, a agents (e.g. nurses), r resources (e.g. rooms), and a time horizon of T units of time, the computational complexity for solving the constraint portion of the formulation (i.e. Equations (2)–(11)) would be of the order of $O\left(2^{\binom{|{}^tA_{\tau_i}^a| + |{}^tR_{\tau_i}^r| + |H_{\tau_i}|}{\tau_i}}\right) = O(2^{(anmT)(rnmT)(n)}) = O(2^{an^3m^2T^2r})$.

In the typical example for the Labor and Delivery Floor we consider in this paper, one might expect a stressful scenario to include approximately $n = 20$ patients each with $m = 3$ stages of labor, $a = 10$ nurses, $r = 20$ rooms, a planning horizon of 12 hours with an event rate of one event every 10 minutes (i.e. $T = 12 \text{ hours} * \frac{1 \text{ event}}{10 \text{ minutes}} * \frac{60 \text{ minutes}}{1 \text{ hour}} = 72$). Under such a scenario, the computational complexity of the search for a solution considering only the constraints, would be $O(2^{10*20^3*3^2*72^2*20}) = O(2^{74,649,600,000})$.

4.1. The role of the resource nurse

The functions of a resource nurse are to assign nurses to take care of labor patients and to assign patients to labor beds, recovery room beds, operating rooms, antepartum ward beds, or postpartum ward beds. The resource nurse has substantial flexibility when assigning beds, and their decisions will depend upon the type of patient and the current status of the unit in question. They must also assign scrub technicians to assist with surgeries in operating rooms, and call in additional nurses if required. The corresponding decision variables for staff assignments and room/ward assignments in the above formulation are ${}^tA_{\tau_i}^a$ and ${}^tR_{\tau_i}^r$, respectively.

The resource nurse may accelerate, delay, or cancel scheduled inductions or cesarean sections in the event that the floor is too busy. Resource nurses may also request expedited active management of a patient in labor. The decision variables for the timing of transitions between the various steps in the care process are described by $s_{\tau_i}^j$ and $f_{\tau_i}^j$. The commitments to a patient (or that patient's procedures) are represented by H_{τ_i} .

The resource nurse may also reassign roles among nurses: for example, a resource nurse may pull a nurse from triage or even care for patients herself if the floor is too busy; or if a patient's condition is particularly acute (e.g. the patient has severe preeclampsia), the resource nurse may assign one-to-one nursing. The level of attentional resources a patient requires and the level a nurse has available correspond to variables $U_{\tau_i}^j$ and ${}^tG_{\tau_i}^a$, respectively. The resource nurse makes their decisions while considering current patient status $\Lambda_{\tau_i}^j$, which is manually transcribed on a whiteboard, as shown in Figure 2.

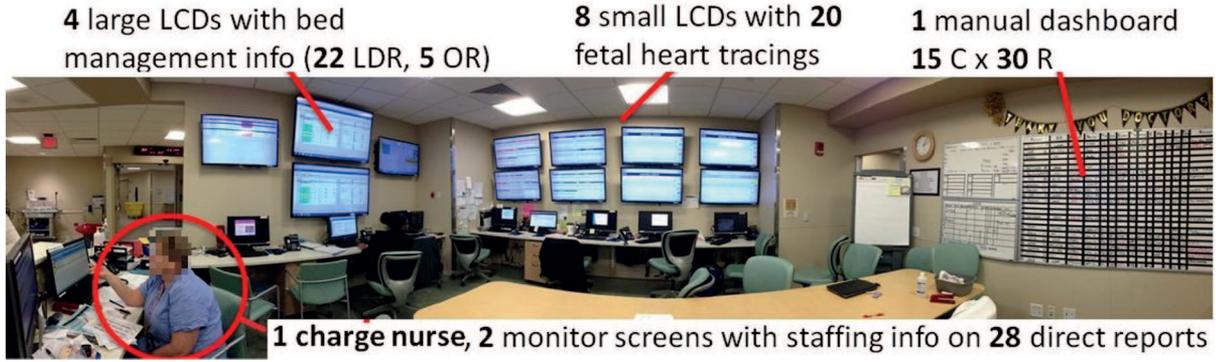


Fig. 2. A resource nurse must assimilate a large variety and volume of information to effectively reason about resource management for patient care.

5. Implementation of decision support

There are two fundamental challenges to providing decision support guidance through direct solution of the optimization problem depicted above. First, the computational complexity of the problem precludes production of real-time solutions. The computational complexity of completely searching for a solution satisfying constraints in Equations (2)–(11) is given by $O\left(2^{|A||R|T^2} C_a^{|A|T}\right)$, where $|A|$ is the number of agents, with each agent possessing an integer-processing capacity of C_a ; there are n tasks τ_i , each with m_i subtasks; $|R|$ resources; and an integer-valued planning horizon of T units of time. In practice, there are ~ 10 nurses (agents) who can care for up to two patients at a time (i.e. $C_a = 2, \forall a \in A$), 20 different rooms (resources) of varying types, 20 patients (tasks) at any one time and a planning horizon of 12 hours or 720 minutes, yielding a worst-case complexity of $\sim 2^{10 \times 20 \times 720^2} 2^{10 \times 720} \geq 2^{10^6}$, which is computationally intractable.

The second challenge to decision support guidance is that the precise form of the objective function (Equation (1)) that resource nurses optimize for is unknown. Prior work has indicated that domain experts are adept at describing the (high-level, contextual, and task-specific) features used in their decision making, yet it is more difficult for experts to describe how they reason about these features (Cheng et al., 2006; Raghavan et al., 2006). As such, we applied a machine learning technique to learn a set of heuristic scheduling policies from demonstrations of resource nurse decision making. We then applied these learned policies to produce advice for the computer-based and robotic decision support systems.

5.1. Learning from resource nurses

In this section, we present a framework for learning (via expert demonstration) a set of heuristics for resource allocation and scheduling that emulates resource nurse decision making. For the purposes of our experiment, we focused on learning a policy for recommending which nurse

should care for which patient, and for making patient room assignments. We demonstrate in our results section that this technique produced high-quality recommendations, as evidenced by an overall 90% accept rate of high-quality advice.

We applied action-driven learning rather than explicitly modeling a reward function and relying upon dynamic programming or constraint solvers. This latter approach (Abbeel and Ng, 2004; Konidaris et al., 2011; Odom and Natarajan, 2015; Vogel et al., 2012; Zheng et al., 2014; Ziebart et al., 2008) can quickly become computationally intractable for problems involving hundreds of tasks and tens of agents due to memory limitations. Approximate dynamic programming approaches exist that essentially reformulate the problem as regression (Konidaris et al., 2011; Mnih et al., 2015), yet the amount of data required to regress over a large state space remains challenging, and Markov decision process (MDP)-based task allocation and scheduling solutions exist only for simple problems (Wu et al., 2011; Wang and Usher, 2005; Zhang and Dietterich, 1995).

Instead, we applied an apprenticeship scheduling algorithm (Gombolay et al., 2016a) inspired by work in Web page ranking (Jin et al., 2008; Pahikkala et al., 2007). The model representation, a graph with nodes and directed arcs, provides a suitable analogy for capturing the complex temporal dependencies (i.e. precedence, wait, and deadline constraints) relating tasks within a scheduling problem. The approach uses pairwise comparisons between the actions taken (e.g. scheduling agent a to complete task τ_i at time t) and the set of actions not taken (e.g. unscheduled tasks at time t) to learn relevant model parameters and scheduling policies demonstrated by the training examples. One key advantage of this pairwise approach is that it is non-parametric, in that the cardinality of the input vector is not dependent upon the number of tasks (or actions) that can be performed in any instance.

Consider a set of *task–resource–agent* $\left\langle \tau_i^j, R_{\tau_i^j}^a, A_{\tau_i^j}^a \right\rangle$ assignments, denoted $\pi_q \in \Pi$. Each assignment π_q has

a set of associated features, γ_{π_q} , indicating patient type (i.e. patients presenting with scheduled induction, scheduled cesarean section, or acute unplanned anomaly), bed type, whether or not the bed is occupied, and staff status (i.e. the number of patients for which the staff member is serving as primary nurse, covering nurse, baby nurse, or scrub technician). Next, consider a set of m observations, $O = \{O_1, O_2, \dots, O_m\}$. Each observation consists of a feature vector describing the task–resource–agent tuple π_q scheduled by the expert demonstrator (including a null task τ_\emptyset , resource r_\emptyset or agent a_\emptyset if no task, resource or agent was scheduled). The goal is to then learn a policy that correctly determines which task–resource–agent tuple π_q to schedule as a function of feature state.

To learn to correctly assign the subsequent task to the appropriate resource and/or agent, we transform each observation O_m into a new set of observations by performing pairwise comparisons between the scheduled assignment π_q and the set of assignments s that were not scheduled:

$$\begin{aligned} \text{rank}_{\langle \pi_q, \pi_r \rangle}^m &:= [\gamma_{\pi_q} - \gamma_{\pi_r}], & y_{\langle \pi_q, \pi_r \rangle}^m &= 1, \\ \forall \pi_r \in \Pi \setminus \pi_q, \forall O_m \in \mathcal{O} | \pi_q \text{ scheduled in } O_m & \end{aligned} \quad (14)$$

$$\begin{aligned} \text{rank}_{\langle \pi_r, \pi_q \rangle}^m &:= [\gamma_{\pi_r} - \gamma_{\pi_q}], & y_{\langle \pi_r, \pi_q \rangle}^m &= 0, \\ \forall \pi_r \in \Pi \setminus \pi_q, \forall O_m \in \mathcal{O} | \pi_q \text{ scheduled in } O_m & \end{aligned} \quad (15)$$

Equation (14) creates a positive example for each observation in which a π_q was scheduled. This example consists of the input feature vector, $\phi_{\langle \pi_q, \pi_r \rangle}^m$, and a positive label, $y_{\langle \pi_q, \pi_r \rangle}^m = 1$. Each element of the input feature vector $\phi_{\langle \pi_q, \pi_r \rangle}^m$ is computed as the difference between the corresponding values in the feature vectors γ_{π_q} and γ_{π_r} , describing scheduled assignment π_q and unscheduled task π_r . Equation (15) creates a set of negative examples with $y_{\langle \pi_r, \pi_q \rangle}^m = 0$. For the input vector, we take the difference of the feature values between unscheduled assignment π_r and scheduled assignment π_q .

We applied these observations to train a decision-tree classifier $f_{\text{priority}}(\pi_q, \pi_r) \in \{0, 1\}$ to predict whether it is better to make the task–resource–agent assignment π_q as the next assignment rather than π_r . Given this pairwise classifier, we can determine which single assignment π_q^* is the highest-priority assignment according to

$$\widehat{\pi_q^*} = \arg \max_{\pi_q \in \Pi} \sum_{\pi_r \in \Pi} f_{\text{priority}}(\pi_q, \pi_r) \quad (16)$$

by determining which assignment is most often of higher priority in comparison with the other assignments in Π .

In our experiments, we directly applied $f_{\text{priority}}(\pi_q, \pi_r)$ to generate high-quality recommendations. We generated low-quality advice using two methods. The first method recommended the action that minimized Equation (16), instead of maximizing it. This approach would typically generate infeasible advice (e.g. to move a patient to a room that

is currently occupied). A second method was applied to offer low-quality but feasible advice (e.g. to assign a post-operating patient to triage). This was achieved by evaluating Equation (16) after filtering the space of possible actions to include only feasible actions (per the constraints in Equations (2)–(11)). As examples, advice given as high-quality, low-quality feasible, and low-quality infeasible could be, respectively, to “assign a new patient to a nurse who currently has no patients under her care;” “assign a new patient to a nurse who already has patients under her care;” and “assign a new patient to an occupied room.” We note that recommendations for the low-quality condition were produced by randomly selecting between the infeasible and feasible methods to mitigate ordering effects.

The dataset used for training was generated by seven resource nurses working with the simulation for a total of $2\frac{1}{2}$ hours, simulating 60 hours of elapsed time on a real labor floor. This yielded a dataset of more than 3,013 individual decisions. None of the seven resource nurses who contributed to the dataset participated in the experiment.

As we note throughout the presentation of the experiment, a high rate of acceptance of advice (e.g. advice given in the “high-quality advice” condition) does not imply the advice is, in fact, high quality. Nonetheless, we do believe there is evidence to support that we present participants with a meaningful gradation of advice quality. First, the algorithm has been demonstrated to learn a high-quality representation of human decision-making in prior work (Gombolay et al., 2016b). Second, we internally validated the learned policy, $f_{\text{priority}}(\pi_q, \pi_r)$, through expert review with a team of obstetricians and labor and delivery professionals. Third, we investigated the various strategies for decision-making of the resource nurses in our dataset in a recent paper (i.e. Molina et al. (2018) published within the medical community. For this paper, we internally validated that the dataset accurately captured the decision-making process present in labor and delivery operations.

6. Results

We report statistical testing of our hypotheses here. We defined statistical significance at the $\alpha = 0.05$ level.

6.1. Analysis and discussion of H1

Objective measures of compliance and reliance were assessed based on the participant’s “accept” or “reject” responses to each decision support recommendation. Statistics for hits, misses, false alarms, and correct rejections are reported in Table 1. We note that the robot- and computer-based decision support systems provided a total of 412 and 417 suggestions, respectively, across 24 participants and two conditions (i.e. high- and low-quality advice) per participant for an average of 8.58 and 8.69 suggestions, respectively, per participant per condition.

Table 1. Confusion matrix for participants shown as a raw count and percentage of the whole.

Robotic Decision support		Response	
		Accept	Reject
Advice quality	High	188 (45.6%)	20 (4.9%)
	Low	26 (6.3%)	178 (43.2%)

Computer Decision support		Response	
		Accept	Reject
Advice quality	High	176 (42.2%)	27 (6.5%)
	Low	21 (5.0%)	193 (46.2%)

Table 2. Rates of correct “accept” and “reject” responses (and the corresponding Type I and Type II error rates), as well as the PPV and NPV for participants, depicted as percentages. Note that the rate of correct “accept” is one minus than the Type I error rate, and the rate of correct “reject” is one minus than the Type II error rate.

	Robot	Computer	p-value
Correct Accept Rate (Type I Error Rate)	90.4% (9.6%)	86.7% (13.3%)	0.241
Correct Reject Rate (Type II Error Rate)	87.3% (12.7%)	90.2% (9.8%)	0.343
PPV	87.9%	89.3%	0.635
NPV	89.9%	87.7%	0.483

As shown in Table 2, results from a z-test for two proportions indicated no statistically significant difference in the Type II error rates between the robotic ($p_R = 12.7\%$) and computer-based ($p_C = 9.8\%$) decision support conditions ($z = 0.949$, $p = 0.343$) nor in the Type I error rates ($p_R = 9.6\%$, $p_C = 13.3\%$, $z = 1.174$, $p = 0.241$). For reference, the corresponding rates of correct “accept” responses to high-quality advice (i.e. 1 minus the Type II error rate) are $p_R = 90.4\%$, $p_C = 86.7\%$, and the corresponding rates of correct “reject” responses to low-quality advice (i.e. 1 minus the Type I error rate) are $p_R = 87.3\%$, $p_C = 90.2\%$. Further, results from a z-test for two proportions indicated no statistically significant difference in positive predictive value (PPV) between the robotic ($p_R = 87.9\%$) and computer-based ($p_C = 89.3\%$) decision support conditions ($z = 0.474$, $p = 0.635$) nor in negative predictive value (NPV) between the robotic ($p_R = 89.9\%$) and computer-based ($p_C = 87.7\%$) decision support conditions ($z = 0.702$, $p = 0.483$). Results from a two one-sided tests (TOST) equivalence test using two z-tests for two proportions indicated that the Type I error rate, PPV, and NPV were statistically equivalent between the robotic and virtual decision support conditions.

We also analyzed the rates of Type I and Type II errors in the second and third trials, upon transition in advice quality (Table 3). Results from a z-test for two proportions indicated a significant difference in the rate of incorrect “accept” of low-quality advice (Type I error) across the second and third trials for the computer-based decision support (6.7% versus 17.6%, $z = 1.793$, $p = 0.036$), but not for the

Table 3. Correct accept and reject decisions made with computer-based (computer-accept, computer-reject) versus robotic (robot-accept, robot-reject) decision support, as a function of trial number, depicted as a raw count and percentage of the whole.

	Trial number			
	Bad advice		Good advice	
	1	2	3	4
Computer-accept	5 (10.4%)	4 (6.7%)	41 (82.0%)	49 (92.5%)
Robot-accept	9 (17.6%)	5 (9.6%)	49 (87.5%)	44 (93.6%)

	Trial number			
	Good advice		Bad advice	
	1	2	3	4
Computer-reject	7 (16.3%)	7 (12.3%)	42 (82.4%)	52 (94.5%)
Robot-reject	8 (14.8%)	2 (3.9%)	42 (85.7%)	47 (90.4%)

robotic support (9.6% versus 16.0%, $p = 0.461$). In other words, participants’ rates of Type I error associated with computer-based support increased significantly when they had received high-quality advice in the previous trial.

Moreover, findings from a contrast test indicated that the average Type I error rates among participants under the computer-based decision support condition after transitioning from high- to low-quality advice ($p_C^{after} = 18.4\%$) were the highest among the set of Type I error rates for both the robotic and computer-based decision support systems across the second and third trials ($p_C^{before} = 9.7\%$, $p_C^{after} = 18.4\%$, $p_R^{before} = 9.7\%$, $p_R^{after} = 14.4\%$, $p = 0.002$). In other words, rates of Type I errors were highest when participants had received high-quality advice during the previous trial and were now receiving low-quality advice from the computer-based decision support system.

Similarly, contrast test results showed that the average Type II error rates among participants under the computer-based decision support condition after transitioning from low- to high-quality advice ($p_C^{after} = 22.6\%$) were highest among the set of Type II error rates for both the robotic and computer-based decision support systems across the second and third trials ($p_C^{before} = 10.6\%$, $p_C^{after} = 22.6\%$, $p_R^{before} = 3.2\%$, $p_R^{after} = 12.0\%$, $p = 0.002$). In other words, rates of Type II errors were highest when participants had received low-quality advice in the previous trial and were now receiving high-quality advice from the computer-based decision support system.

6.2. Analysis and discussion of H2

A composite measure of trust was computed, as in the work by Jian et al. (2000). Results from a repeated-measures analysis of variance (RANOVA) demonstrated a statistically significant increase in the average rating for the decision support system under the high-quality advice condition ($M = 5.17$) as compared with the low-quality condition

Table 4. Subjective measures: post-trial questionnaire items for which participants' responses were statistically significantly different across experimental conditions. Responses to Questions 1–5 were on a 7-point (reverse) scale; Questions 4–10 were on a 10-point scale.

Trust and Embodiment in Human–Robot Interaction

1. The system is deceptive.
 2. I am suspicious of the system's intent, action, or outputs.
 3. The system provides security.
 4. Distant/close.
 5. I think this decision support coach could be a friend of mine.
 6. I think I could have a good time with this decision support coach.
 7. People will find it interesting to use this decision support coach.
 8. People will find this decision support coach attractive.
 9. While you were interacting with this decision support coach, how much did you feel as if it were a social being?
 10. Unsociable/sociable.
-

($M = 3.22$), with a standard error estimate of $SD = 9.551$, ($F(1, 23) = 82.848, p < 0.001$).

However, a RANOVA yielded no statistically significant difference in user trust between the robotic ($M = 4.22$) and computer-based ($M = 4.16$) embodiment conditions, with a standard error estimate of $SD = 0.264$, ($F(1, 23) = 0.070, p = 0.707$). Results from a TOST equivalence test, using two t -tests, indicated that subjects' trust ratings for the computer-based and robotic support were within one point of each other on a seven-point Likert scale.

We observed significant differences in the attitudinal assessment of the robotic versus computer-based decision support conditions for Questions 4–10 in Table 4, indicating that participants rated the robotic system more favorably. The result was established using a two-way omnibus Friedman test, followed by pairwise Friedman tests. The test statistics for the pairwise Friedman tests were $p = 0.033$, $p = 0.048$, $p = 0.038$, $p < 0.001$, $p = 0.046$, $p = 0.023$, and $p = 0.002$, respectively. Strikingly, there was not a single question (out of 37) for which participants rated the computer-based decision support significantly better than the robotic support.

We also found that the subjective perception of the character of the robot was significantly less sensitive to transitions in advice quality than the computer-based decision support. We computed the frequency with which the ratings of one embodiment condition subsumed the other, and vice versa. Specifically, we defined $x_{R,L}$ as the Likert-scale rating for a given question and a particular participant under the robotic low-quality advice condition, and likewise for the high-quality condition, $x_{R,H}$. The variables $x_{C,L}$ and $x_{C,H}$ were similarly defined for the computer-based low- and high-quality conditions. The robotic condition was defined as subsuming the computer-based condition if either $\min(x_{R,L}, x_{R,H}) \leq \min(x_{C,L}, x_{C,H}) \leq$

$\max(x_{C,L}, x_{C,H}) < \max(x_{R,L}, x_{R,H})$ or $\min(x_{R,L}, x_{R,H}) < \min(x_{C,L}, x_{C,H}) \leq \max(x_{C,L}, x_{C,H}) \leq \max(x_{R,L}, x_{R,H})$, and vice versa for the computer-based condition subsuming the robotic condition. A χ^2 test indicated that the participants' subjective evaluation according to Questions 1, 2, and 3 ($p = 0.047, 0.046$, and 0.047 , respectively) changed more significantly under the computer-based condition than the robotic condition. There were no questions for which the response changed more significantly under the robotic condition versus the computer-based condition. In other words, the subjective assessment of the robot was more robust to advice quality changes than the computer-based decision support. Further investigation is warranted to determine whether these effects persist over time as users habituate to interaction with the robot.

6.3. Post-hoc analysis of reaction times

In post-hoc analysis, we investigated the effects of the mode of decision support and the quality of advice given by the decision support system on participants' reaction time to that advice. We conducted a two-factor, RANOVA to test whether participants took more time responding to advice as a function of advice quality (i.e. high versus low quality) and the mode of support (i.e. robotic versus computer-based).

The result of the RANOVA suggested a relationship ($F(1, 720) = 3.15, p = 0.077$) between advice quality and participants' response times. Specifically, participants took 2.351 ± 2.280 seconds when responding to low-quality advice versus 2.048 ± 2.148 seconds when responding to high-quality advice, corresponding to a 15% increase in time taken to respond to low-quality advice.

The evidence provided by the RANOVA was less compelling for the relationship between the mode of decision support and participants' response times ($F(1, 720) = 0.278, p = 0.598$): participants took 2.135 ± 2.065 seconds when responding to computer-based decision support versus 2.267 ± 2.148 seconds when responding to the robotic system, corresponding to a 19% increase in time taken. However, a possible relationship between operator response time and embodiment is an area for future work.

6.4. Discussion

The statistical analysis provided in the results section show that the Type I and Type II error rates were comparable between robotic and computer-based decision support systems. Furthermore, embodiment appeared to offer performance gains, as indicated by lower error rates after the quality of recommendation changed mid-experiment. These encouraging findings provide evidence that a robotic assistant may be able to participate in decision making with nurses without eliciting inappropriate dependence. One potential rationale for these results is that experts may be less susceptible to the negative effects of embodiment, as

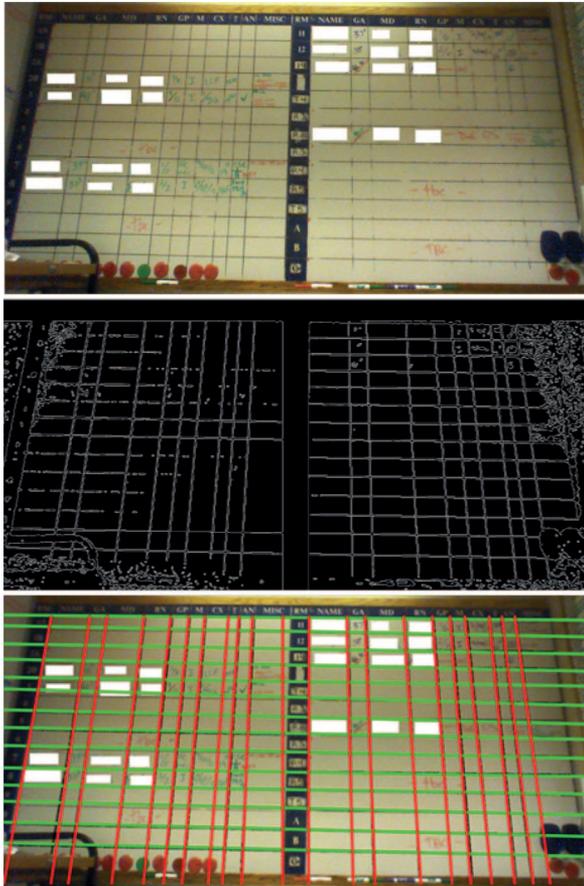


Fig. 3. (Top) An anonymized photo of a labor floor dashboard, taken by the camera from a Nao robot. (Middle) The same image after the preprocessing and edge detection steps are applied. (Bottom) The segmentation recovered after the Hough transform and postprocessing steps are applied.

has been documented previously among experienced users interacting with anthropomorphic agents (Pak et al., 2012). In addition, note that our study was conducted with a stationary robot, with movement limited to co-speech gestures. Further investigation is warranted for situations in which experts interact with mobile service robots that participate in decision making.

The statistical analysis for H1 focuses on determining whether there are differences in the Type I and Type II error rates between the robot- and computer-based decision support systems. We found statistically significant differences between these error rates when a condition of low- or high-quality advice was followed by the opposite condition. We found that the transition point results in an inflated Type I or II error rate for participants working with the computer-based decision support system as opposed to the robot-based decision support system. It is notable that this statistically significant difference arises through change in the embodiment of the decision support system, and that participants working with the computer-based decision support experience higher rates of inappropriate compliance

and reliance when the advice quality changes. The possible presence of even a transient inflation of error rates in a safety-critical domain, such as a labor and delivery ward, is potentially significant and warrants careful evaluation. Future studies will be required to assess whether this effect is transient or whether participants re-calibrate to the new advice quality over time.

Our findings support H2 in that the robotic system was rated more favorably on attitudinal assessment than computer-based decision support, even as it engendered appropriate dependence. It is inevitable that a service robot will occasionally make poor-quality suggestions, and we positively note that the robot engendered greater tolerance of errors than the computer-based decision support. These results indicate a positive signal for successful adoption of a robot that participates in a resource nurse's decision making.

In our post-hoc analysis studying the effect of advice quality and embodiment and reaction time, the results were mixed. We found positive evidence that reaction times were degraded when participants were evaluating low-quality advice. These data are supported by prior work in neuroscience by Goodyear et al. (2016) examining advice response time as a function of advice quality. On the other hand, although the average reaction time for considering advice from the robot-based decision support system were delayed as opposed to the computer-based decision support system, the statistical analysis was not significant. This inconclusive result was also found in prior work by Goodyear et al. (2016). As such, we have established a new hypothesis, to be tested in a future experiment. First, we hypothesize that operator response time is negatively proportional to advice quality. Second, embodied (i.e. robotic) decision support elicits more attention, as measured by a longer response time, compared with un-embodied (i.e. computer-based) support.

There is also a noteworthy correspondence between our observation that the robot engendered greater tolerance of errors than a computer-based decision support and prior work in shared decision-making authority in human-robot teaming (Gombolay et al., 2015) culminating from a series of prior experimental investigations (Gombolay et al., 2014; Gombolay and Shah, 2014a,b; Gombolay et al., 2013). In our experiment, we found that participants' subjective perception of a robotic decision support system varied less than that of a computer-based system as a function of the quality of advice offered by that system. In related work, Gombolay et al. (2015) investigated a scenario in which a human would work on a manufacturing team with the participant, a second human teammate (i.e. a confederate), and a third agent, either another human or a robot. The team would work under one of three conditions: (1) the participant would allocate tasks to the team (i.e. determine which team member would complete which tasks); (2) the third agent would allocate tasks to the team; or (3) the participant and the third agent would share allocation responsibility. Gombolay et al.

(2015) found that participants' subjective perception of the value of the third agent varied more significantly when that agent was a robot than a human. In other words, the human agent engendered a greater level of tolerance for the task allocation schema than the robotic agent. In this paper, we report that a robotic agent engendered greater tolerance for decision support errors than did a computer-based system. Although the factors in these experiments differ, there is interesting evidence to support the idea that anthropomorphizing affects the tolerance or variability of an operator's perception of that agent. As such, we establish the following hypothesis for future work: As an agent becomes increasingly anthropomorphic, a participant's perception of that agent grows increasingly tolerant of its behavior or role in the interaction.

7. Pilot demonstration of a robotic assistant on the labor and delivery floor

Based on the positive results of our experiment, we conducted a pilot demonstration in which a robot-assisted resource nurses on a labor and delivery floor at a tertiary care center.

7.1. Robot system architecture

The system comprised three primary subsystems providing the vision, communication, and decision support capabilities, respectively. Figure 4 depicts a system diagram.

7.1.1. Vision system. In our experiments, the statuses of patients, nurses, and beds were provided and updated within the simulation. In contrast, nurses and robots on a real labor floor must read handwritten information off of a whiteboard (i.e. "dashboard"), as depicted in Figures 2 and 3. Extracting and parsing this information autonomously with high degrees of accuracy and reliability presents a substantial technical challenge. We made two assumptions to address this: (1) that the set of physician and nurse names is closed and known in advance, and (2) that patient names are transcribed for the robot upon patient arrival.

In our demonstration, we leveraged the structured nature of the dashboard to introduce priors that ensured patient information was interpretable. Rows on the dashboard indicate room assignments, while columns indicate patient parameters (e.g. attending physician, gestational age, etc.). Once our robot captured an image of the dashboard on the labor and delivery floor, we applied a Canny edge detection operator (Canny, 1986) and Hough transformation (Duda and Hart, 1972) to isolate the handwriting in individual grid cells (Figure 3). To improve segmentation accuracy, images were preprocessed to remove colored writing, the header row at the top of the whiteboard, and the room name columns at the left and middle of the whiteboard. The images were also postprocessed using heuristics to improve the location of the line estimates. Specifically, rows were

biased toward having uniform heights, whereas columns of each field were biased to be the same width on the left and right side of the central column that contained room numbers.

The contents of each grid cell were processed using a classification technique appropriate to the data type therein. Numeric fields were parsed using a convolutional neural network (CNN)² trained on MNIST data, whereas alphabetical fields with known sets of possible values (e.g. attending physician, nurse names) were parsed using a multi-class CNN trained on handwriting (Figure 5).

Handwriting samples consisting of 28 uniquely written alphabets served as a basis for generating classifier training data. Fonts were created from the provided samples and used (along with system fonts) to create a large set of binary images containing samples of nurse names. These synthetic writing samples were constructed with a range of applied translations, scalings, and kerning values within a 75×30 pixel area.

The vision system was used to determine the current status of patient–nurse allocations, nurse role information, and room usage. Prior to deployment, we performed a validation of the vision system and found our recognition system to correctly classify handwritten samples across 15 classes (names) with $\sim 83.7\%$ overall accuracy and 97.8% average accuracy. These results were obtained without performing any environmental manipulations, such as adjusting lighting or employing high-resolution cameras. In the pilot deployment, our vision system assisted humans with transcription of patient data by presenting suggested input to the user for confirmation.

7.1.2. Communication. The system separated auditory processing (i.e. speech recognition) into two components before determining an appropriate query and issuing it to the decision support system: the first component transcribed a user's spoken commands to text; the second converted the text into an appropriate, object-based JAVA query to be issued to the apprenticeship scheduler to generate a decision support recommendation.

To transcribe a user's spoken commands to text, we employed CMUSphinx,³ an open-source speech-to-text recognition software developed by Carnegie Mellon University. To achieve high-level performance in a live setting, we defined a list of template-based phrases a user might utter, such as "Where should I move the patient in room [#]?" or "Who should nurse [Name] take care of?" To understand the person's query, the system merely needed to determine (1) which type of query the person had made and (2) which keywords define the parameters of that query (i.e. the room number in the query, "Where should I move the patient in room [#]?").

Based on information available a priori (e.g. the list of nurse names), the system enumerated all possible instantiations of these template-based queries and added them to the CMUSphinx's language database. We modified this to more

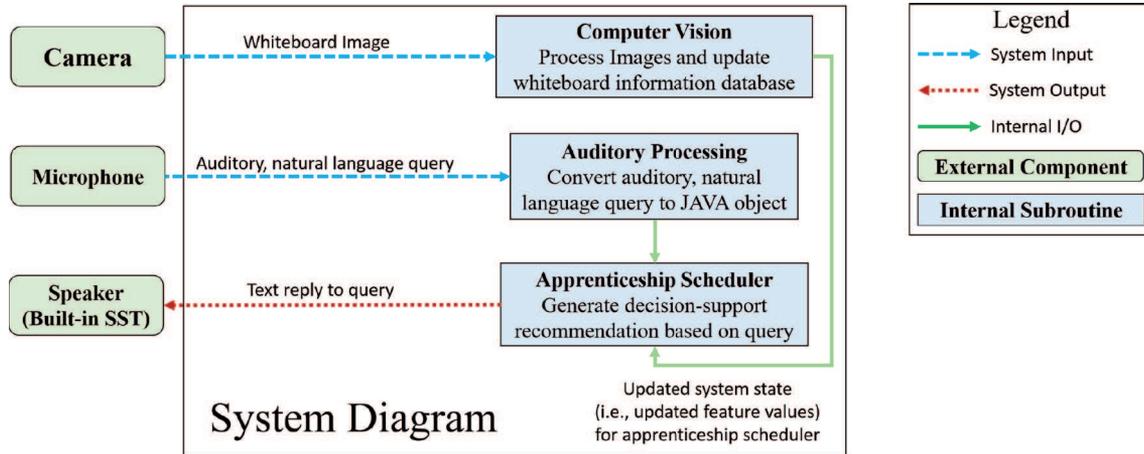


Fig. 4. The architecture of the robotic decision support system.

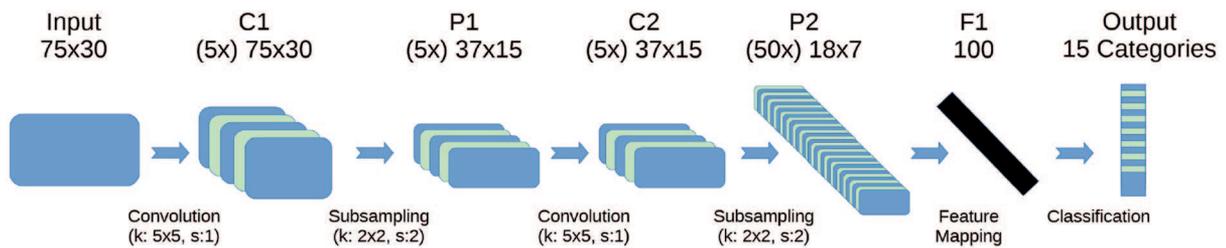


Fig. 5. The convolutional neural network architecture we utilized for our vision system. This network comprised a single input layer for an individual grid cell from the white board, two alternating sets of convolutional (C1 and C2) and maxpool layers (P1 and P2), followed by a single fully connected layer (F1), with 15 output classes: one for each of the nurses' names.



Fig. 6. The robot system in action on the labor floor.

heavily weight the likelihood of words that were part of the aforementioned key phrases to increase the probability of proper transcription. We also required a supplementary pronunciation file to compile the database for runtime.

When a user needed to query the robotic decision support system, they would press a button on a laptop to initiate the auditory processing subroutine. The laptop transmitted the audio from the microphone to the processing subroutine until that subroutine determined that the user had finished their query. The CMUSphinx module would listen to the audio in real time and constantly update its prediction as the user spoke: upon completion, the system produced a text-based transcription of the query.

To determine which query type was spoken by the user, the system removed all keywords defining the parameters of the query from the transcript (e.g. the room number from the example above), then computed the Levenshtein distance (Levenshtein, 1966) between the pruned transcript and each of the query templates. The system then rank-ordered the query templates based on their Levenshtein distance from the pruned transcript, and selected the one with the highest rank as the inferred query, with one caveat: if the highest ranking query template's keyword options (e.g. a room number) did not match the transcribed keyword (e.g. the name of a nurse), the system eliminated the highest-ranking query, then considered the next-highest

query. This process would repeat until the first satisfiable query template was identified.

Finally, based on this inference, the auditory processing subroutine would call the apprenticeship scheduler, providing as input the type of query and an ordered list of the keywords that parameterize the query.

7.1.3. Decision support. The live pilot demonstration of the robot incorporated the same mechanism for generating decision support as that used during our experiments. However, unlike the experiments, the decision support system's input was taken from the vision subsystem, and the user query from the communication subsystem. The set of possible actions to be recommended was filtered according to the query as recognized by the communication subsystem. For example, if the user asked "Where should I move the patient in room 1A?," actions that would change nurse assignments were not considered. The recommended action was communicated to the user via text-to-speech software.

7.2. Feedback from nurses and physicians

We conducted a test demonstration on the labor floor (Figure 6). Three users interacted with the robot over the course of 3 hours. Ten queries were posed to the robot; seven resulted in successful exchanges and three failed due to background noise. A live recording of the demo is available at <http://tiny.cc/RobotDemo>.

After interacting with the robotic support, User 1, a physician, said, "I think the [robot] would allow for a more even dispersion of the workload amongst the nurses. In some hospitals ...more junior nurses were given the next patient...more senior nurses were allowed to only have one patient as opposed to two." User 2, a resource nurse, said, "New nurses may not understand the constraints and complexities of the role, and I think the robot could help give her an algorithm ...that she can practice, repeat, and become familiar with so that it becomes second nature to her." User 3, a labor nurse, offered, "I think you could use this robot as an educational tool."

8. Future work

The experiment results and successful system demonstration raise several promising areas of future work. First, further study is needed to understand how embodiment affects trust, reliance, and compliance for a spectrum of decision support form factors. In this experiment, two form factors were investigated: a computer-based decision support system and a relatively anthropomorphic robotic system (i.e. an Aldebaran NAO). Based upon our initial findings, one may hypothesize that a more humanoid robot (e.g. a Honda ASIMO) might further improve the ability of participants to detect changes in advice quality.

Second, future work will investigate how the relationship between the nursing staff and the decision support system

evolves over days and weeks as a part of a longitudinal study. The cross-sectional study we report in this paper provides initial insight into the effects of embodiment on trust, reliance, and compliance. However, it is important to evaluate how these phenomena might evolve with time. Such a longitudinal study also presents technical challenges. As the robot's experience grows, so too should the robustness of its machine learning representation of the ideal patient flow through the hospital. Future research will improve upon our machine learning formulation to continuously incorporate new experiences in a semi-supervised fashion so as to minimize the need for tedious human labeling while not sacrificing the quality of the policy.

Finally, future research is necessary to design the robotic decision support system to function as a training tool. Perhaps by learning from experts, the robot could serve as an infinitely patient tutor for novice nurses or other support staff. We found initial evidence in prior work regarding the efficacy of our apprenticeship scheduling algorithm serving as a tutor in a military defense training simulation (Gombolay et al., 2017). In future work, we aim to investigate the efficacy of such a tutor for managing patient flow in hospitals.

9. Conclusion

This paper addresses two barriers to fielding intelligent hospital service robots that take initiative to participate with nurses in decision making. We observed experimental evidence that experts performing decision-making tasks may be less susceptible to the negative effects of support embodiment. Further, our decision support was able to produce context-specific decision strategies and apply them to make reasonable suggestions for which tasks to perform and when. Finally, based on the previous two findings, we conducted a first successful test demonstration in which a robot assisted resource nurses on a labor and delivery floor in a tertiary care center.

Author(s) note

Matthew Gombolay is currently affiliated to Georgia Institute of Technology, Atlanta, GA, USA.

Funding

This work was supported by the National Science Foundation Graduate Research Fellowship Program under grant number 23883577, CRICO Harvard Risk Management Foundation, and Aldebaran Robotics Inc.

Notes

1. The Kaiser Family Foundation has provided an interactive database of employment statistics at: <http://kff.org/other/state-indicator/total-number-of-professionally-active-nurses-by-gender/?currentTimeframe=0>.

2. Thanks go to Mikhail Sironenko for developing this package, which is available at <https://sites.google.com/site/mihailsironenko/projects/cuda-cnn>.
3. CMU Sphinx Open Source Speech Recognition Toolkit; available at <http://cmusphinx.sourceforge.net/>.

References

- Abbeel P and Ng AY (2004) Apprenticeship learning via inverse reinforcement learning. In: *ICML*. New York: ACM Press.
- Bainbridge WA, Hart JW, Kim ES and Scassellati B (2011) The benefits of interactions with physically present robots over video-displayed agents. *International Journal of Social Robotics* 3(1): 41–52.
- Bartneck C, Kuli'c D, Croft E and Zoghbi S (2009) Measurement instruments for the anthropomorphism, animacy, likeability, perceived intelligence, and perceived safety of robots. *International Journal of Social Robotics* 1(1): 71–81.
- Bertsimas D and Weismantel R (2005) *Optimization over Integers*. Belmont: Dynamic Ideas.
- Bloss R (2011) Mobile hospital robots cure numerous logistic needs. *Industrial Robot: An International Journal* 38(6): 567–571.
- Brandenburg L, Gabow P, Steele G, Toussaint J and Tyson BJ (2015) Innovation and best practices in health care scheduling. Technical report, Institute of Medicine of the National Academies.
- Canny J (1986) A computational approach to edge detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 8(6): 679–698.
- Chen JY, Barnes MJ and Harper-Sciarini M (2011) Supervisory control of multiple robots: Human-performance issues and user-interface design. *IEEE Transactions on Systems, Man, and Cybernetics, Part C: Applications and Reviews* 41(4): 435–454.
- Cheng TH, Wei CP and Tseng VS (2006) Feature selection for medical data mining: Comparisons of expert judgment and automatic approaches. In: *Proceedings of CBMS*, pp. 165–170.
- Cummings ML and Guerlain S (2007) Developing operator capacity estimates for supervisory control of autonomous vehicles. *Human Factors: The Journal of the Human Factors and Ergonomics Society* 49(1): 1–15.
- de Visser EJ, Krueger F, McKnight P, et al. (2012) The world is not enough: Trust in cognitive agents. *Proceedings of the Human Factors and Ergonomics Society Annual Meeting* 56(1): 263–267.
- Desai M, Kaniarasu P, Medvedev M, Steinfeld A and Yanco H (2013) Impact of robot failures and feedback on real-time trust. In: *Proceedings of the 8th ACM/IEEE International Conference on Human-Robot Interaction (HRI '13)*. Piscataway, NJ: IEEE Press, pp. 251–258.
- Desai M, Medvedev M, V'azquez M, et al. (2012) Effects of changing reliability on trust of robot systems. In: *2012 7th ACM/IEEE International Conference on Human-Robot Interaction (HRI)*. Piscataway, NJ: IEEE Press, pp. 73–80.
- DiGiuse N (2013) Hospitals hiring robots. Available at: http://www.electronicproducts.com/Computer_Peripherals/Systems/Hospitals_hiring_robots.aspx.
- Dismukes RK, Berman BA and Loukopoulous LD (2007) *The Limits of Expertise: Rethinking Pilot Error and the Causes of Airline Accidents*. Ashgate Publishing.
- Dixon SR and Wickens CD (2006) Automation reliability in unmanned aerial vehicle control: A reliance–compliance model of automation dependence in high workload. *Human Factors: The Journal of the Human Factors and Ergonomics Society* 48(3): 474–486.
- Duda RO and Hart PE (1972) Use of the hough transformation to detect lines and curves in pictures. *Communications of the ACM* 15(1): 11–15.
- Gombolay M, Gutierrez R, Clarke S, Sturla G and Shah J (2015) Decision-making authority, team efficiency and human worker satisfaction in mixed human-robot teams. *Autonomous Robots* 39(3): 293–312.
- Gombolay M, Gutierrez R, Sturla G and Shah J (2014) Decision-making authority, team efficiency and human worker satisfaction in mixed human-robot teams. In: *Proceedings of Robots: Science and Systems (RSS)*, Berkeley, CA.
- Gombolay M, Jensen R, Stigile J, Son SH and Shah J (2016a) Apprenticeship scheduling: Learning to schedule from human experts. In: *Proceedings of the International Joint Conference on Artificial Intelligence*, New York City, USA.
- Gombolay M, Jensen R, Stigile J, Son SH and Shah J (2017) Learning to tutor from expert demonstration via apprenticeship scheduling. In: *Proceedings of the Association for the Advancement of Artificial Intelligence (AAAI) Workshop on Human–Machine Collaborative Learning (HMCL)*, San Francisco, California.
- Gombolay M and Shah J (2014a) Challenges in collaborative scheduling of human–robot teams. In: *Proceedings of AAAI Fall Symposium Series on Artificial Intelligence for Human–Robot Interaction (AI-HRI)*, Arlington, VA.
- Gombolay M and Shah J (2014b) Increasing the adoption of autonomous robotic teammates in collaborative manufacturing. In: *Proceedings of the Human–Robot Interaction Pioneers Workshop*, Bielefeld, Germany, pp. 62–63.
- Gombolay M, Wilcox R, Artiles AD, Yu F and Shah J (2013) Towards successful coordination of human and robotic work using automated scheduling tools: An initial pilot study. In: *Proceedings of Robots: Science and Systems (RSS) Human–Robot Collaboration Workshop*, Berlin, Germany.
- Gombolay M, Yang XJ, Hayes B, et al. (2016b) Robotic assistance in coordination of patient care. In: *Proceedings of Robotics: Science and Systems (RSS)*, Ann Arbor, MI.
- Goodyear K, Parasuraman R, Chernyak S, Madhavan P, Deshpande G and Krueger F (2016) Advice taking from humans and machines: An fMRI and effective connectivity study. *Frontiers in Human Neuroscience* 10: 542.
- Hu J, Edsinger A, Lim YJ, et al. (2011) An advanced medical robotic system augmenting healthcare capabilities-robotic nursing assistant. In: *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*. Piscataway, NJ: IEEE Press, pp. 6264–6269.
- Jansen N, Cubuktepe M and Topcu U (2017) Synthesis of shared control protocols with provable safety and performance guarantees. In: *ACC*. Piscataway, NJ: IEEE Press, pp. 1866–1873.
- Jian JY, Bisantz AM and Drury CG (2000) Foundations for an empirically determined scale of trust in automated systems. *International Journal of Cognitive Ergonomics* 4(1): 53–71.
- Jin R, Valizadegan H and Li H (2008) Ranking refinement and its application to information retrieval. In: *Proceedings of the 17th International Conference on World Wide Web*. New York: ACM Press, pp. 397–406.

- Kaber DB and Endsley MR (1997) Out-of-the-loop performance problems and the use of intermediate levels of automation for improved control system functioning and safety. *Process Safety Progress* 16(3): 126–131.
- Kehle SM, Greer N, Rutks I and Wilt T (2011) Interventions to improve veterans' access to care: A systematic review of the literature. *Journal of General Internal Medicine* 26(2): 689–696.
- Kidd CD and Breazeal C (2004) Effect of a robot on user perceptions. In: *Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, vol. 4. Piscataway, NJ: IEEE Press, pp. 3559–3564.
- Kiesler S, Powers A, Fussell SR and Torrey C (2008) Anthropomorphic interactions with a robot and robot-like agent. *Social Cognition* 26(2): 169–181.
- Konidaris G, Osentoski S and Thomas P (2011) Value function approximation in reinforcement learning using the Fourier basis. In: *Proceedings of AAAI*. pp. 380–385.
- Kulić D, Venture G, Yamane K, Demircan E, Mizuuchi I and Mombaur K (2016) Anthropomorphic movement analysis and synthesis: A survey of methods and applications. *IEEE Transactions on Robotics* 32(4): 776–795.
- Lee JD and See KA (2004) Trust in automation: Designing for appropriate reliance. *Human Factors* 46(1): 50–80.
- Lee KW, Peng W, Jin SA and Yan C (2006) Can robots manifest personality? An empirical test of personality recognition, social responses, and social presence in human–robot interaction. *Journal of Communication* 56(4): 754–772.
- Levenshtein VI (1966) Binary codes capable of correcting deletions, insertions, and reversals. *Soviet physics doklady* 10(8): 707–710.
- Leyzberg D, Spaulding S and Scassellati B (2014) Personalizing robot tutors to individuals' learning differences. In: *Proceedings of the 2014 ACM/IEEE international conference on Human–Robot interaction*. New York: ACM Press, pp. 423–430.
- Li W, Sadigh D, Sastry SS and Seshia SA (2014) Synthesis for human-in-the-loop control systems. In: *International Conference on Tools and Algorithms for the Construction and Analysis of Systems*. New York: Springer, pp. 470–484.
- Mnih V, Kavukcuoglu K, Silver D, et al. (2015) Human-level control through deep reinforcement learning. *Nature* 518(7540): 529–533.
- Molina RL, Gombolay M, Jonas J, et al. (2018) Association between labor and delivery census and delays in patient management: Findings from a computer simulation module. *Journal of Obstetrics and Gynecology*, in press.
- Murai R, Sakai T, Kawano H, et al. (2012a) A novel visible light communication system for enhanced control of autonomous delivery robots in a hospital. In: *2012 IEEE/SICE International Symposium on System Integration (SII)*. Piscataway, NJ: IEEE Press, pp. 510–516.
- Murai R, Sakai T, Kitano Y and Honda Y (2012b) Recognition of 3D dynamic environments for mobile robot by selective memory intake and release of data from 2D sensors. In: *2012 IEEE/SICE International Symposium on System Integration (SII)*. Piscataway, NJ: IEEE Press, pp. 621–628.
- Mutlu B and Forlizzi J (2008) Robots in organizations: The role of workflow, social, and environmental factors in human–robot interaction. In: *2008 3rd ACM/IEEE International Conference on Human–Robot Interaction (HRI)*. Piscataway, NJ: IEEE Press, pp. 287–294.
- Odom P and Natarajan S (2015) Active advice seeking for inverse reinforcement learning. In: *Proceedings of AAAI*, pp. 4186–4187.
- Olsen DR Jr and Wood SB (2004) Fan-out: measuring human control of multiple robots. In: *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*. New York: ACM Press, pp. 231–238.
- Özkil AG, Fan Z, Dawids S, Aanes H, Kristensen JK and Christensen KH (2009) Service robots for hospitals: A case study of transportation tasks in a hospital. In: *IEEE International Conference on Automation and Logistics, 2009 (ICAL'09)*. Piscataway, NJ: IEEE Press, pp. 289–294.
- Pahikkala T, Tsivtsivadze E, Airola A, Boberg J and Salakoski T (2007) Learning to rank with pairwise regularized least-squares. In: *SIGIR 2007 Workshop on Learning to Rank for Information Retrieval*, pp. 27–33.
- Pak R, Fink N, Price M, Bass B and Sturre L (2012) Decision support aids with anthropomorphic characteristics influence trust and performance in younger and older adults. *Ergonomics* 55(9): 1059–1072.
- Parasuraman R, Sheridan T and Wickens CD (2000) A model for types and levels of human interaction with automation. *IEEE Transactions on Systems, Man and Cybernetics, Part A: Systems and Humans* 30(3): 286–297.
- Pizer SD and Prentice JC (2011) What are the consequences of waiting for health care in the veteran population? *Journal of General Internal Medicine* 26(2): 676–682.
- Raghavan H, Madani O and Jones R (2006) Active learning with feedback on features and instances. *Journal of Machine Learning Research* 7: 1655–1686.
- Robinette P, Howard AM and Wagner AR (2016) Overtrust of robots in emergency evacuation scenarios. In: *Proceedings of HRI*.
- Schröder M and Trouvain J (2003) The German text-to-speech synthesis system mary: A tool for research, development and teaching. *International Journal of Speech Technology* 6(4): 365–377.
- Sheridan TB (2011) Adaptive automation, level of automation, allocation authority, supervisory control, and adaptive control: Distinctions and modes of adaptation. *IEEE Transactions on Systems, Man and Cybernetics, Part A: Systems and Humans* 41(4): 662–667.
- Shipman SA and Sinsky CA (2013) Expanding primary care capacity by reducing waste and improving efficiency of care. *Health Affairs (Millwood)* 32(11): 1990–1997.
- Takayama L and Pantofaru C (2009) Influences on proxemic behaviors in human-robot interaction. In: *IEEE/RSJ International Conference on Intelligent Robots and Systems, 2009 (IROS 2009)*. Piscataway, NJ: IEEE Press, pp. 5495–5502.
- Tapus A, Tapus C and Mataric M (2009) The role of physical embodiment of a therapist robot for individuals with cognitive impairments. In: *Proceedings of the International Symposium on Robot and Human Interactive Communication (RO-MAN)*. Piscataway, NJ: IEEE Press, pp. 103–107.
- Vogel A, Ramach D, Gupta R and Raux A (2012) Improving hybrid vehicle fuel efficiency using inverse reinforcement learning. In: *Proceedings of AAAI*, pp. 384–390.

- Wang YC and Usher JM (2005) Application of reinforcement learning for agent-based production scheduling. *Engineering Applications of Artificial Intelligence* 18(1): 73–82.
- Wickens CD, Hollands JG, Banbury S and Parasuraman R (2013) *Engineering psychology and human performance*. New York: Pearson Education.
- Wickens CD, Li H, Santamaria A, Sebok A and Sarter NB (2010) Stages and levels of automation: An integrated meta-analysis. *Proceedings of the Human Factors and Ergonomics Society Annual Meeting* 54(4): 389–393.
- Winston WL, Venkataramanan M and Goldberg JB (2003) *Introduction to mathematical programming*, vol. 1. Pacific Grove, CA: Thomson/Brooks/Cole Duxbury.
- Wu J, Xu X, Zhang P and Liu C (2011) A novel multi-agent reinforcement learning approach for job scheduling in grid computing. *Future Generation Computer Systems* 27(5): 430–439.
- Zhang W and Dietterich TG (1995) A reinforcement learning approach to job-shop scheduling. In: *Proceedings of the International Joint Conference on Artificial Intelligence*, pp. 1114–1120.
- Zheng J, Liu S and Ni L (2014) Robust Bayesian inverse reinforcement learning with sparse behavior noise. In: *Proceedings of AAAI*, pp. 2198–2205.
- Ziebart BD, Maas A, Bagnell JA and Dey AK (2008) Maximum entropy inverse reinforcement learning. In: *Proceedings of AAAI*, pp. 1433–1438.

Appendix. Index to multimedia extensions

Archives of IJRR multimedia extensions published prior to 2014 can be found at <http://www.ijrr.org>, after 2014 all videos are available on the IJRR YouTube channel at <http://www.youtube.com/user/ijrrmultimedia>

Table of Multimedia Extension

Extension	Media type	Description
1	Video	Nao offering advice to a participant with speech and co-speech gestures.
2	Video	Tutorial describing the labor and delivery floor simulation.