

Balancing the Mental Load: Adaptive Human-Agent Approaches for Peak Performance

Deborah van Sinttruije
Delft University of Technology
Delft, Netherlands
D.vanSinttruije@tudelft.nl

Frans A. Oliehoek
Delft University of Technology
Delft, South Holland, Netherlands
f.a.oliehoek@tudelft.nl

Catharine Oertel
Delft University of Technology
Delft, South Holland, Netherlands
C.R.M.M.Oertel@tudelft.nl

ABSTRACT

In high-stakes environments, where errors have severe consequences, designing adaptive systems that adjust in real-time to a user's cognitive state is valuable but challenging due to the non-observable nature of these states. This study explores a partially observable Markov decision process (POMDP) framework to infer hidden cognitive states and dynamically manage cognitive load, minimizing the risk of cognitive overload. We tested two models: one based on established literature and another fine-tuned with user data using the Bayesian Particle Marginal Metropolis-Hastings (PMMH) method. Both models were evaluated against a performance-adaptive baseline in a user study. Our findings show that POMDP-based agents significantly reduce errors, improve task performance over time, and provide a more balanced perceived task difficulty. These results suggest that while POMDP-based adaptive systems can improve human performance, future work on cognitive adaptive systems should focus on refining model estimation techniques to better capture individual cognitive states.

KEYWORDS

Cognitive Load, Working Memory, POMDP

1 INTRODUCTION

Human-agent collaboration has gained increasing attention, driven by research into how AI models can complement human capabilities [2, 7, 31, 32]. This synergy is promising, particularly because humans face cognitive limitations. For instance, in high-stakes tasks such as emergency responses, AI-supported decision systems can assist operators by presenting task-relevant information in stressful environments, thereby reducing workload during time-critical procedures [36]. Yet, if the AI delivers excessive or overly complex information, the human operator risks cognitive overload, potentially leading to fatal errors.

Cognitive load, the mental effort required to process and retain information, is closely linked to the capacity of working memory [53], which serves as a temporary storage system managing information relevant to the task at hand. However, working memory has a limited capacity that constrains the amount of information that can be actively maintained at any given time [42]. When cognitive load surpasses this momentary capacity, task performance diminishes, resulting in errors and reduced efficiency [53]. This is particularly problematic in high-stakes environments, where maintaining situational awareness without overwhelming cognitive resources is essential [25, 49]. Therefore, the development of

cognitive adaptive systems could enhance performance by accurately assessing and responding to the user's current cognitive capabilities, thereby preventing overload and enhancing overall effectiveness [22, 41].

Significant progress has been made toward developing cognitive adaptive systems, with recent studies predicting cognitive load associated with specific tasks [12, 38]. However, real-time systems that can estimate an individual's cognitive load and adapt accordingly remain challenging to implement. This difficulty stems from the fact that internal cognitive states are not directly observable and vary based on individual and contextual factors [23]. For instance, cognitive load is influenced not only by task complexity but also by emotional states like affect, as well as fatigue, stress, and task familiarity [11, 13, 23, 29]. While previous attempts to model these cognitive states exist [17, 20, 27, 39], conventional methods that adapt by increasing task complexity until failure often fail to capture the nuances of individual cognitive capacity.

In this paper, we explore an adaptive framework where an AI agent adjusts working memory task difficulty based on the user's cognitive capacity. The framework is tested in a simulated high-stress environment with varying stress levels to assess its robustness. The AI dynamically adjusts task difficulty to maintain optimal cognitive load and performance. Since cognitive states are not directly observable, we use a partially observable Markov decision process (POMDP) to infer hidden cognitive states based on user actions and feedback. The POMDP's ability to gather information over time refines cognitive load estimates, enhancing task adaptation. Through reinforcement learning, the agent determines an optimal policy to manage cognitive load and minimize overload risk.

We implement two models: one based on established cognitive capacity theories and a second fine-tuned using task-specific data through the Bayesian Particle Marginal Metropolis-Hastings (PMMH) method, proven effective in hidden Markov models (HMMs) [46]. The effectiveness of these models is evaluated through a user study comparing them to a standard adaptive algorithm, measuring both objective performance and participants' perceived task difficulty. This comparison demonstrates the benefits of dynamic cognitive load management. Our hypotheses include:

- H1 The POMDP framework reduces user errors compared to the standard adaptive method, while maintaining task difficulty.
- H2 Users perceive tasks as more balanced in difficulty when using the POMDP framework compared to the traditional method.
- H3 The POMDP framework achieves results in less error over time in comparison to the traditional method.

- H4 The fine-tuned model leads to better performance outcomes and a more positive user experience compared to the theory-based model.

Our findings show that the POMDP model offers significant advantages over traditional methods in memory tasks, reducing user errors while balancing actual and perceived task difficulty. Although the fine-tuned model outperformed traditional methods, it did not significantly surpass the theory-based model. These results highlight the POMDP framework’s potential, but further research is needed to enhance its effectiveness with real user data.

This paper contributes the following:

(1) Adaptive Model Design: A novel working memory model that predicts an individual’s capacity and enables real-time dynamic adjustments. (2) Bayesian Fine-Tuning: A sophisticated Bayesian method to fine-tune the POMDP model with user-specific data, potentially improving cognitive load predictions. (3) Evaluation Against Traditional Methods: A user study comparing the POMDP agent to conventional methods, showing potential to reduce errors and balance task difficulty, with areas for further improvement.

2 BACKGROUND AND RELATED WORK

This section reviews cognitive adaptive systems and POMDPs in Human-Computer Interaction (HCI), with a focus on how cognitive load management and adaptive mechanisms improve task performance and user experience. The first part examines cognitive adaptation strategies, while the second explores the application of POMDPs in modeling complex human behaviors under uncertainty.

2.1 Cognitive Adaptive Systems

Early applications of cognitive adaptive systems appeared in high-stakes environments [22, 25, 26], focusing on optimizing dynamic interactions between human operators and machine interfaces under critical conditions. These systems typically followed a generalized, one-size-fits-all approach, adapting to operators’ mental workload, environmental factors, and mission criteria due to the lack of reliable methods to measure individual cognitive differences [47]. While these early systems established key principles for adaptive interface design, they often lacked the sophistication and personalization needed to account for cognitive variability.

Subsequent research has made significant progress in modeling and detecting cognitive load on an individual basis [5, 8, 35, 44, 47]. Furthermore, studies across various domains demonstrate the benefits of integrating cognitive load assessments into adaptive systems. For instance, Chan et al. [15] developed a cognitive context-aware system that improved user receptivity to memory tasks under low cognitive load, while Lindlbauer et al. [41] adapted mixed reality designs based on users’ cognitive states. Predictive models, such as Kelleher and Hnin [38], enhance adaptive learning tools by forecasting cognitive load, and Barral et al. [8] demonstrated the feasibility of real-time user modeling.

Despite these advancements, a gap remains in research that both develops user models and implements corresponding interventions [58]. Many systems focus on detecting user states and predicting challenges but fall short in dynamically adjusting the task environment or workload. While the groundwork for cognitive adaptive

systems is in place, the challenge of inferring hidden cognitive states and adapting to them in real time remains unresolved.

2.2 POMDP models in HCI

POMDPs are particularly well-suited for capturing complex human behaviors, such as those influenced by mental states or intentions, due to their ability to model latent states [16, 30, 52, 54, 55]. As a result, they provide a natural solution for inferring hidden cognitive states. Moreover, their ability to manage uncertainty and adapt dynamically to human interactions has led to widespread use in HCI [34] and human-robot collaboration (HRC) [40], particularly for modeling complex decision-making under uncertain conditions. A key advantage of POMDPs is their capacity for active data gathering, enabling the model to make informed decisions based on current observations while considering actions that could enhance future outcomes. This ability to gather and process information in real time makes POMDPs an intriguing tool for designing adaptive, user-centered systems that can respond to evolving cognitive states.

Despite the benefits of the POMDP framework, applying it to real-world problems remains challenging due to the complexity of estimating latent model parameters from literature and the difficulty of deriving the true POMDP model from data [56]. Traditional methods, such as the Baum-Welch algorithm, rely on maximum likelihood estimation, which often fails to adequately account for parameter uncertainty [19].

In response, Bayesian methods have gained traction, offering a probabilistic framework that handles uncertainty in model parameters and incorporates prior knowledge. This makes them particularly useful in complex scenarios, providing a more robust and accurate alternative for POMDP model estimation. Recent research has demonstrated the potential of Bayesian methods in the estimation of cognitive models [37], highlighting their potential in related areas. Among these techniques, Particle Marginal Metropolis-Hastings (PMMH) [4] has emerged as a powerful method for parameter estimation in Hidden Markov Models (HMMs) [19]. PMMH, a Bayesian Markov Chain Monte Carlo (MCMC) algorithm, combines the strengths of particle filtering with MCMC’s robustness, allowing efficient sampling from the posterior distribution of model parameters.

Given the success of PMMH in HMMs, it is reasonable to hypothesize that this method could also be effective for POMDPs. HMMs are similar to POMDPs in that both involve latent states and uncertainty in observations, but HMMs lack the action component present in POMDPs. This similarity suggests that PMMH could provide a more robust and scalable alternative to the traditional Baum-Welch approach for POMDPs. By leveraging the probabilistic nature of Bayesian inference, PMMH can explore the parameter space more thoroughly and provide more accurate estimates, especially in scenarios where model complexity or large data sets present significant challenges.

3 METHODS

In this section, we describe the development, implementation, and evaluation of two POMDP models for managing cognitive load in high-stakes tasks, compared against a baseline in a user study (see Figure 1 for an overview of the models).

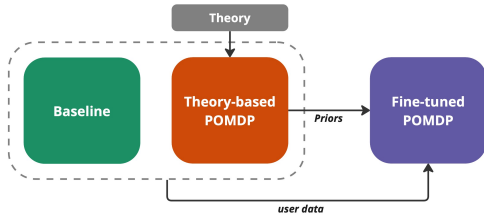


Figure 1: The setup of the agent conditions in the study. We define a baseline and theory-based POMDP model. The fine-tuned model is based on the user data from both the baseline condition and the theory based model.

This study centers around a wire-cutting task, inspired by collaborative problem-solving in bomb defusing scenarios and the game Keep Talking and Nobody Explodes [28]. A visual representation of the task is shown in Figure 2. In the task, an AI agent communicates a sequence of wires to the participant, who then has to remember the sequence and solve a puzzle by cutting the correct wires. Participants are then tested for their recollection of the wires. This process was repeated twelve times within a 30-minute timeframe to measure cognitive fatigue without overburdening participants [1]. To simulate varying levels of stress, the task becomes more demanding midway through the trials.

At each time step, the agent determines how much information to relay, balancing the risk of cognitive overload with the need for sufficient information. The models were evaluated on their ability to manage this balance, ensuring participants could process the information effectively without becoming overwhelmed. The following sections detail the development of agent conditions and the design of the user study.

3.1 Baseline

In the baseline condition, the agent was equipped with a performance-adaptive algorithm for working memory, as described by Woods et al. [60]. This algorithm dynamically adjusts task difficulty based on the participant’s performance to better match their working memory capacity. When the participant responds correctly, the difficulty increases, and the agent adds one more wire to remember. If the participant fails twice in a row, the agent reduces the number of wires by one.

3.2 Theory-based POMDP

To establish a robust model for predicting human performance in our wire-cutting task, we developed a theory-driven POMDP model based on established research in cognitive psychology and human factors. Specifically, the model incorporates key factors that influence cognitive load, including task familiarity, fatigue, and stress.

Task Familiarity. Research shows that as individuals become more familiar with a task, their cognitive load decreases because routine processes become more automated, allowing them to handle

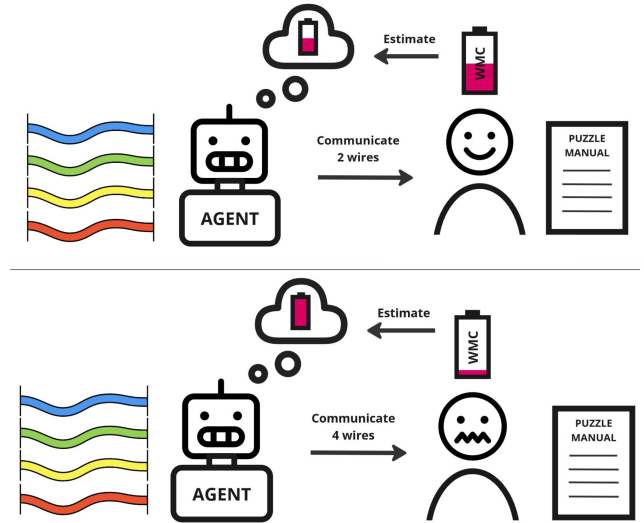


Figure 2: The wire cutting task. At the top half we see an agent that correctly estimates the users current working memory capacity. At the bottom half, the agent overestimates the users working memory capacity and relays too many wires.

more complex aspects without overloading working memory [29, 51, 57]. In our model, we assume that participants’ working memory capacity increases as they progress through the task, reflecting their growing familiarity with its demands.

Fatigue. Mental fatigue affects cognitive load management over time, reducing the mental resources available for processing information and learning [11, 48]. Our model accounts for this by simulating a gradual decline in working memory capacity as a function of time, reflecting the accumulation of cognitive fatigue participants may experience during prolonged tasks.

Stress. Stress can indirectly increase cognitive load by impairing focus and narrowing attention [14], making individuals more susceptible to distractions, which often leads to poorer performance [10]. Additionally, cognitive resources may be diverted to managing emotional responses rather than concentrating on the task [13, 23]. Our model reflects this by making a decline in working memory capacity more likely under stressful conditions, simulating the impact of stress on cognitive performance.

We formally define our working memory POMDP as follows:

Definition 3.1 (Working Memory POMDP). The Working Memory (WM) POMDP is tuple $(S, A, \Omega, T, O, R, \gamma)$, where:

- S is the set of states $s: \langle WMC, tr \rangle$. There are 8 working memory capacities $WMC = 3, \dots, 10$ and 12 trials $tr = 1, \dots, 12$. The WMC represents the number of items the user can remember, starting with 3 items, as defined in the working memory performance adaptive algorithm [60]. The maximum number of items is 10, which was the highest working memory capacity reached in a pilot study ($n=31$). The 12

trials represent the total number of trials the user must complete. At each step, the environment is in a state $s \in S$.

- A is the set of actions $\{3, \dots, 10\}$. There are 8 different actions, each corresponding to one of the 8 WMC states. These actions, taken by the agent, represent the number of items the agent assigns for the user to remember.
- Ω is the set of observations $\{\text{incorrect}, \text{correct}\}$. These represent the agent’s possible observations after assigning a number of items to the user. The user’s response is either correct or incorrect.
- T is the transition matrix $|S| \times |A| \times |S|$, where $T[s', a, s] = p(s'|a, s)$, which is the probability of transitioning to state s' , given that the agent takes action a in state s . Transition probabilities were modeled using a Gaussian distribution to simulate progression over time steps. The initial state distribution b_0 is assumed to be uniform.
- O is the observation matrix $|\Omega| \times |A| \times |S|$, where $O[o, a, s] = p(o|a, s)$, thus the probability of observation o , given action a in state s . There was a fixed 10% probability of receiving an incorrect observation.
- R is the reward function for the state-action pair $|S| \times |A|$. The agent receives a reward of +100 when its assigned number of items matches the user’s working memory capacity. If the agent assigns more items than the user can handle, it receives a penalty of -50, whereas assigning fewer items incurs a smaller penalty of -10. Rewards are provided at each transition and accumulated over the course of an episode.
- γ is discount factor of 0.95. The discount factor determines how future rewards are weighted relative to immediate rewards.

In the 12 trials, we simulated the effects of task familiarity, stress, and fatigue by skewing the transition distribution. During time steps 1-5, working memory capacity was more likely to expand as the user became familiar with the task, so the transition distributions skewed toward higher capacity states. At time step 6, a stressor reduced cognitive resources, leading to a decrease in working memory capacity. Similarly, at time steps 9-12, working memory capacity was more likely to decline due to fatigue, with the transitions skewed toward lower capacity states.

The agent’s policy was designed to maximize user performance by dynamically adjusting task difficulty in response to inferred cognitive load, ensuring the task remained challenging but not overwhelming. The inferred cognitive load, or the agent’s belief, was represented as a probability distribution over possible working memory capacity states, as shown in Figure 3.

3.3 Fine-tuned POMDP

To fine-tune our model, we used data from both the baseline and working memory POMDP conditions. Model parameters were estimated using the PMMH algorithm with Adaptive Metropolis covariance adaptation, as outlined by [59]. The priors for the parameters were based on our theory-based model, with a standard deviation of 0.2, prioritizing estimates that deviated by less than 20% prior model to aid model conversion.

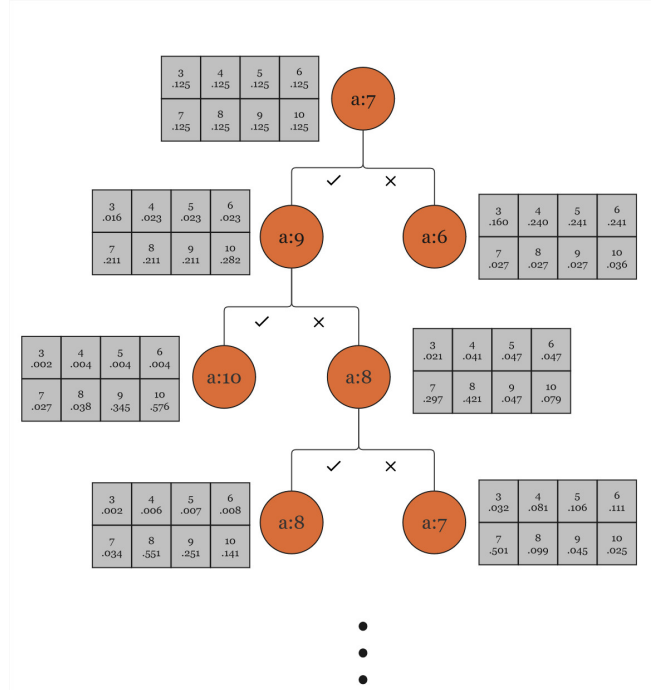


Figure 3: (Partial) Working Memory POMDP policy. In the gray matrices, we see the belief of the agents represented by the possible working memory capacity states, and accompanying probabilities. The nodes represent the actions that the agent takes (the number of wires the agent assigns to the user). The user can either respond correct (check) or incorrect (cross). Given this response, the agent updates its belief for the next step.

3.4 User Evaluation

We conducted a user study on Prolific to evaluate our model. A total of 120 participants were recruited, with 40 participants assigned to each condition: (1) the baseline, (2) the theory-based POMDP, and (3) the fine-tuned POMDP. The sample was balanced by gender (50% male, 50% female), covered a wide age range (18–50), and included individuals from diverse cultural backgrounds while ensuring adequate English proficiency. The study followed ethical guidelines and was approved by our institutional Research Ethics Board.

3.5 Evaluation Metrics

The following metrics were used to evaluate task performance, cognitive load, and user experience across the experimental conditions. A summary of the variables is provided in Table 1.

Task Performance. Task performance was measured by the accuracy with which participants recalled the wires. Accuracy for each trial was recorded as 0 for incorrect recollection and 1 for correct recollection.

Variable	Description
Agent Condition	The agent: baseline, theory-based POMDP, or fine-tuned POMDP.
Trial Number	The specific trial number (1-12) in the experiment.
Stress Condition	Two phases: lower stress in trials 1-6 and higher stress trials 7-12.
Task Performance	Accuracy of wire recall (0 = incorrect, 1 = correct).
Task Difficulty	Number of wires assigned; more wires indicate greater difficulty.
Perceived Difficulty	Subjective rating of difficulty: "Easy," "Moderate," or "Difficult."
Stress-Level Ratings	Stress rated per stress condition on a 0-10 visual analog scale (VAS).

Table 1: Summary of Evaluation Metrics and Variables

Task Difficulty. Objective task difficulty was determined by the number of wires assigned to the participant, with more wires indicating higher difficulty. The agent tailored this difficulty to match the user’s working memory capacity.

Perceived Difficulty. Participants rated the difficulty of remembering the wires on a 3-point scale: "Easy," "Moderate," and "Difficult."

Stress-Level Ratings. Stress levels were measured using a visual analog scale (VAS), with 0 indicating "no stress" and 10 representing "the highest imaginable level of stress." This method, adapted from Barre et al. [9], has been validated for measuring emotional states, including stress.

3.6 Study Procedure

After providing informed consent, participants began the experiment by completing a brief trial to familiarize themselves with the wire-cutting task. After which they completed the first block of six trials, where the agent communicated information for participants to decide which wire to cut. After each trial, neutral feedback was provided ("Correct" or "Incorrect").

Following the first block, pressure was introduced using a modified design from Sattizahn et al. [50]. Participants were informed that they had been paired with another participant and could win \$3 if they both completed the next block faster while maintaining accuracy. The pressure was further increased in the second block of six trials by displaying a red timer during each trial and providing more intense feedback ("OK" for correct and "WRONG" in red for incorrect), following the method of Almazrouei et al. [3].

After each block, participants’ perceived stress was measured using a Visual Analog Scale (VAS). After each trial perceived task difficulty was measured.

3.7 Analysis Procedure

Overall Difficulty and Performance. The overall objective *performance* and objective *difficulty* across the three *agent* conditions were compared. Since the data were not normally distributed, group means were compared using the Kruskal-Wallis test. This analysis addresses the first and partially the fourth research question.

Perceived Difficulty. The overall *perceived difficulty* across all agent groups was compared using the Chi-squared (χ^2) test. This analysis contributes to answering the second research question and partially addresses the fourth research question.

Performance Over Time. To explore the effect of the *agent* on *performance* over *time*, we will visually inspect agent trajectories and conduct a logistic regression analysis. Additionally, we will examine how *stress* responses over *time* might be mitigated by the agents using a repeated measures ANOVA. These analyses help answer the third research question and contribute to the partial answer of the fourth research question.

4 RESULTS

Overall performance. Figure 4 shows the overall *performance*, i.e., the sum of correct answers per participant. A Kruskal-Wallis test revealed a significant difference between the *agent* groups ($p < 0.05$). Post hoc pairwise comparisons, conducted using Dunn’s test with Bonferroni correction, indicated that both the theory-based model and the fine-tuned model resulted in significantly higher accuracy compared to the Baseline ($p < 0.05$ for both theory-based Model vs. Baseline and fine-tuned Model vs. Baseline). No significant difference in accuracy was found between the theory-based Model and the fine-tuned Model.

Additionally, no significant differences were found in the mean number of wires provided by the agents i.e. overall *difficulty* across the agent groups, indicating that the POMDP agents did not consistently assign easier or more difficult tasks compared to the Baseline condition. The mean number of wires per participant for each group are shown in Figure 5.

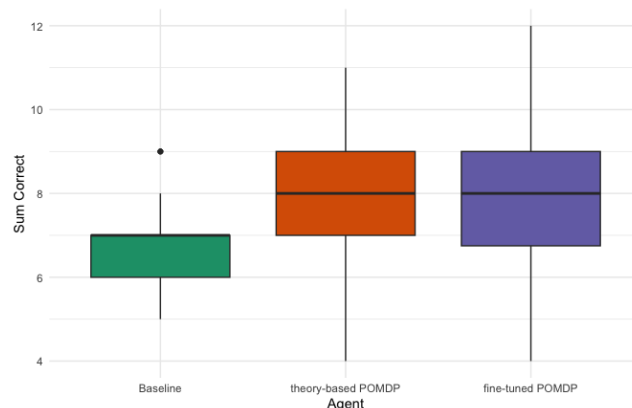


Figure 4: The sum of correct responses of each participant per agent condition.

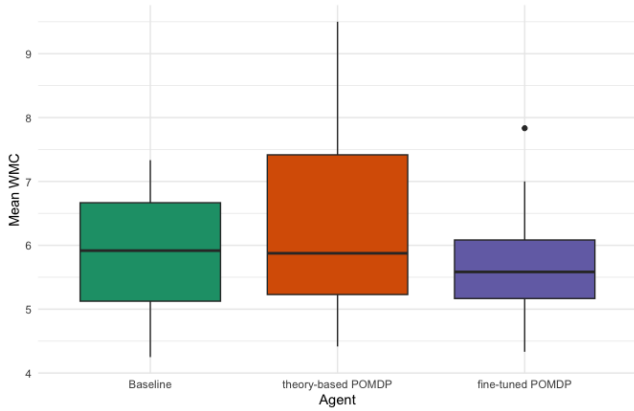


Figure 5: The mean difficulty i.e. the mean number of wires over all trials for each participant per agent condition.

Perceived difficulty. A Chi-square test of independence was conducted to examine the distribution of *perceived difficulty* levels (Difficult, Moderate, Easy) across the *agent* groups, revealing a significant difference, $\chi^2(4, N = 3) = 16.76, p < 0.05$. The contingency table is shown in Table 2. Post-hoc pairwise comparisons with Bonferroni correction indicated a significant difference between the Baseline and theory-based POMDP conditions ($p < 0.05$). No significant differences were found between the Baseline and fine-tuned POMDP conditions ($p = 0.742$) or between the theory-based and fine-tuned conditions ($p = 0.160$).

Difficulty	Baseline	Theory-Based	Fine-Tuned	Total
Difficult	319	282	308	909
Moderate	133	185	152	470
Easy	28	13	20	61
Total	480	480	480	1440

Table 2: Comparison of Difficulty Levels Across the three agent conditions

Agent trajectories. The mean number of items provided by the agent per trial, along with the summed correct responses per trial, is shown in Figure 6. The logistic regression analysis revealed several significant findings, detailed in Table 3. *Trial* was found to have a significant negative effect on the *performance* ($p < .05$), indicating that as the trial number increases, the log-odds of correct responses decreases. Both *POMDP* conditions also showed a significant reduction in the log-odds of correct responses compared to the baseline ($p < .05$). However, the interaction between trial and both *POMDP* agents had a significant positive effect ($p < .05$), suggesting that initially, the *POMDP* agents had lower log-odds of correct responses than the baseline. Over time, as the trial number increased, the log-odds of correct responses for the *POMDP* conditions also increased, indicating a higher probability of correct responses as participants progressed through the task.

Parameter	Log-Odds	SE	95% CI	p
(Intercept)	1.26	0.21	[0.85, 1.68]	< .001
trial	-0.16	0.03	[-0.21, -0.10]	< .001
Agent [POMDP]	-1.80	0.29	[-2.39, -1.23]	< .001
Agent [POMDPf]	-1.03	0.29	[-1.61, -0.46]	< .001
trial×Agent [POMDP]	0.34	0.04	[0.26, 0.42]	< .001
trial×Agent [POMDPf]	0.21	0.04	[0.13, 0.29]	< .001

Table 3: Logistic Regression Results

Effect	DFn	DFd	F	p	η^2 (gen)
Agent	2	117	0.4973	0.6095	0.006975
StressCondition	1	117	45.3828	< .001	0.063094
Agent:StressCondition	2	117	11.8928	< .001	0.034092

Table 4: ANOVA Results

Stress. A repeated measures ANOVA examined the effects of *Agent*, *StressCondition*, and their interaction (*Agent:StressCondition*) on *Stress*. The ANOVA table can be found in Table 4. The main effect of *Agent* was not statistically significant. In contrast, a significant main effect of *StressCondition* was found ($p < .05$). Furthermore, the interaction between *Agent* and *StressCondition* was also statistically significant ($p < .05$). This suggests that while the type of *Agent* alone did not significantly influence the outcome, however, that the impact of *StressCondition* varied depending on the level of *Agent*. This interaction effect can be seen in Figure 7.

5 DISCUSSION

In this study, we explored the design and evaluation of an AI agent that adapts task difficulty based on the user’s working memory capacity. We employed a Partially Observable Markov Decision Process (POMDP) framework to model the dynamic interaction between the agent and the user, with the goal of minimizing cognitive overload and improving task performance. The effectiveness of the POMDP agents, both theory-based and fine-tuned, was evaluated by comparing their performance against a baseline adaptive algorithm in a user study.

The POMDP-based agent was specifically designed to operate within the user’s cognitive limits, dynamically adjusting task difficulty to balance cognitive load. Our results show that participants made significantly fewer errors overall when interacting with a POMDP-based agent compared to the baseline algorithm. Notably, the objective difficulty levels assigned by the POMDP agents were not significantly different from those assigned by the baseline. This supports our first hypothesis (H1) and indicates that the agents did not rely on simplistic strategies to reduce errors but instead effectively managed cognitive load to optimize task performance. In addition to objective performance metrics, we examined participants’ perceptions of task difficulty and stress. In the theory-based POMDP condition, users rated the task difficulty as significantly different from the baseline, with more trials rated as “moderate difficulty” rather than “difficult” or “easy,” suggesting better alignment with their working memory capacity. This supports hypothesis

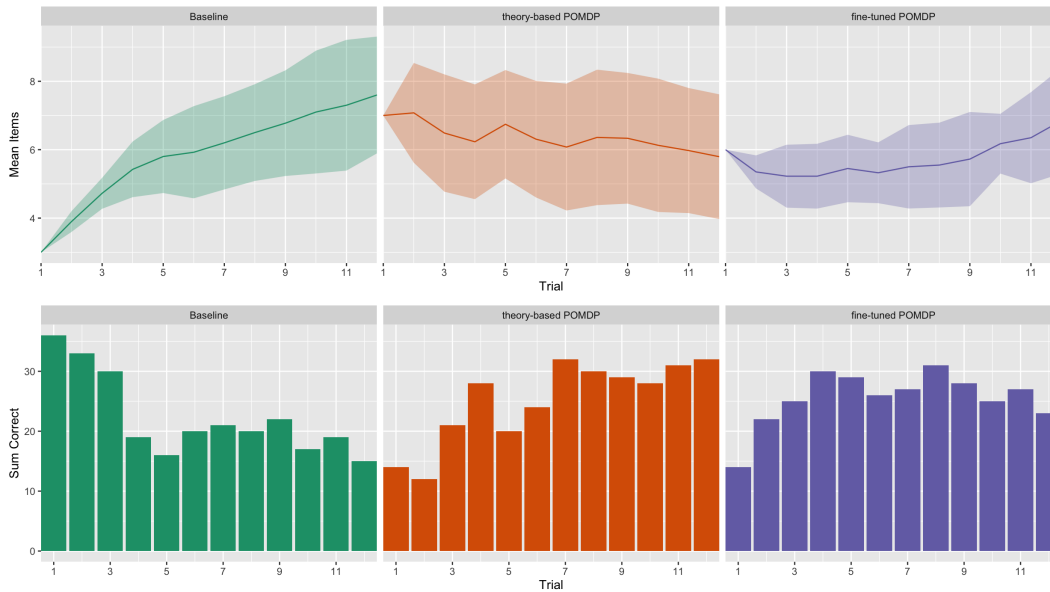


Figure 6: Mean and SD of objective difficulty per trial per agent condition, together with the amount of correct responses per trial.

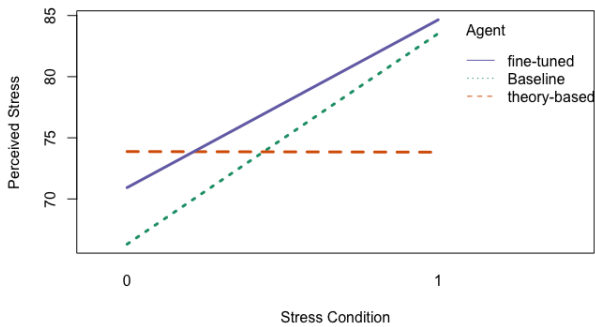


Figure 7: The interaction-effect of stress condition and agent condition on perceived stress.

H2, indicating that the theory-based POMDP agent achieved a more balanced difficulty level. Furthermore, the users’ perception of balanced difficulty suggests that the theory-based agent not only optimized task difficulty in measurable terms but also created a subjective experience that felt appropriately challenging. This aligns with studies on adaptive learning systems, which emphasize the importance of aligning system behavior with user states to improve satisfaction and outcomes [18].

In addition to overall performance metrics, we examined the effect of the adaptive POMDP agents over time. A reduction in errors observed across both POMDP conditions supports our third hypothesis (H3), demonstrating how the POMDP framework actively estimated and adjusted to users’ working memory capacity

in real-time. This reduction in errors, driven by effective cognitive load management, aligns with cognitive load theory [53] and underscores the advantages of cognitive adaptive automation in complex tasks [21]. Moreover, participants in the theory-based POMDP condition appeared less affected by stressful situations, suggesting that the agent’s adaptive adjustments helped mitigate the impact of stress. This reduction in stress, which may have allowed participants to allocate more cognitive resources to the task, potentially boosted performance [13, 23]. These findings align with research showing that well-designed adaptive systems can reduce perceived stress and cognitive load [29, 45].

These findings have important implications for the development of adaptive systems, particularly in areas where managing cognitive load is critical. In addition to high-stakes scenarios, such systems could be valuable in educational technologies, user interfaces, and conversational agents. The ability of POMDP-based systems to dynamically adjust task difficulty based on real-time assessments of working memory capacity suggests that similar approaches could enhance the design of conversational agents. For example, conversational agents could adjust the complexity or amount of information presented to optimize user engagement and comprehension without overwhelming the user [24, 33]. This is especially relevant in educational and assistive technologies, where balancing information delivery with user capacity is crucial for maintaining both effectiveness and user satisfaction [6, 18, 43].

Interestingly, we did not observe significant performance differences between the POMDP agent and the agent fitted with user-specific data. Therefore we cannot accept our final hypothesis (H4). This lack of difference could be attributed to several factors. First, the theory-based agent may have been broad enough to accommodate the variation between participants, making further fine-tuning

unnecessary for this task. Additionally, the fine-tuned agent may not have been optimally fitted or might have been overfitted, resulting in a model that failed to fully capture the nuances of individual cognitive states in new users. These challenges highlight the difficulties in translating advanced algorithms from simulations to real-world applications. In simulation environments, it is both feasible and often recommended to restart the model-fitting process to refine the models [19], as parameters can be iteratively adjusted and validated against simulated outcomes. However, in real-world settings, this iterative validation process is much more challenging. Since the critical variables in these models are often hidden, there is no straightforward way to predict in advance whether a fitted model will perform well in practice. This underscores the need for more robust and thoroughly tested model-fitting methods that can better generalize to real-world conditions.

Despite these challenges, both the theory-based POMDP agent and the fine-tuned POMDP agent significantly outperformed the baseline adaptive algorithm. This suggests that even a basic model within the POMDP framework can offer substantial benefits in managing cognitive load and enhancing user performance, aligning with research that AI systems can augment human performance [2]. The success of the POMDP approach highlights its potential for developing adaptive systems in real-world applications. Future research could explore integrating these models into conversational agents, enabling such systems to not only respond to user inputs but also dynamically anticipate and manage cognitive load, improving both user satisfaction and system effectiveness.

6 CONCLUSION

In this study, we explored the design and implementation of a POMDP framework for adapting to cognitive load through a user study. Participants working with a POMDP-equipped agent demonstrated better accuracy and improved performance over time in a working memory task. Additionally, the theory-based POMDP agent mitigated the effects of stress and provided a more balanced perceived difficulty. However, the fine-tuned POMDP agent, which was adjusted using user-specific data, did not yield additional benefits, suggesting that the fine-tuning process did not significantly improve performance.

These findings highlight the advantages of the POMDP framework in managing cognitive load and enhancing task performance, while also underscoring the challenges of model estimation in real-world applications. Future research should focus on developing robust estimation techniques to better handle the complexities of real-world data.

ACKNOWLEDGMENTS

This research was (partially) funded by the Hybrid Intelligence Center, a 10-year program funded by the Dutch Ministry of Education, Culture and Science through the Netherlands Organisation for Scientific Research.

REFERENCES

- [1] Blair Aitken and Clare MacMahon. 2019. Shared demands between cognitive and physical tasks may drive negative effects of fatigue: A focused review. *Frontiers in Sports and Active Living* 1 (2019), 45.

- [2] Z. Akata, D. Balliet, M. de Rijke, F. Dignum, V. Dignum, G. Eiben, A. Fokkens, D. Grossi, K. Hindriks, H. Hoos, H. Hung, C. Jonker, C. Monz, M. Neerincx, F. Oliehoek, H. Prakken, S. Schlobach, L. van der Gaag, F. van Harmelen, H. van Hoof, B. van Riemsdijk, A. van Wynsberghe, R. Verbrugge, B. Verheij, P. Vossen, and M. Welling. 2020. A research agenda for hybrid intelligence: augmenting human intellect with collaborative, adaptive, responsible, and explainable artificial intelligence. *Computer* 53, 8 (2020), 18–28.
- [3] Mohammed A Almazrouei, Ruth M Morgan, and Itiel E Dror. 2023. A method to induce stress in human subjects in online research environments. *Behavior research methods* 55, 5 (2023), 2575–2582.
- [4] Christophe Andrieu, Arnaud Doucet, and Roman Holenstein. 2010. Particle markov chain monte carlo methods. *Journal of the Royal Statistical Society Series B: Statistical Methodology* 72, 3 (2010), 269–342.
- [5] Tobias Appel, Natalia Sevchenko, Franz Wortha, Katerina Tsarava, Korbinian Moeller, Manuel Ninaus, Enkelejda Kasneci, and Peter Gerjets. 2019. Predicting cognitive load in an emergency simulation based on behavioral and physiological measures. In *2019 International Conference on Multimodal Interaction*. 154–163.
- [6] Ibtissam Azzi, Adil Jeghal, Abdelhay Radouane, and Hamid Tairi. 2019. Personalized e learning systems based on automatic approach. In *2019 International Conference on Wireless Technologies, Embedded and Intelligent Systems (WITS)*. IEEE, 1–6.
- [7] Gagan Bansal, Tongshuang Wu, Joyce Zhou, Raymond Fok, Besmira Nushi, Ece Kamar, Marco Tulio Ribeiro, and Daniel Weld. 2021. Does the whole exceed its parts? the effect of ai explanations on complementary team performance. In *Proceedings of the 2021 CHI conference on human factors in computing systems*. 1–16.
- [8] Oswald Barral, Sébastien Lallé, Grigori Guz, Alireza Iranpour, and Cristina Conati. 2020. Eye-tracking to predict user cognitive abilities and performance for user-adaptive narrative visualizations. In *Proceedings of the 2020 international conference on multimodal interaction*. 163–173.
- [9] Ronan Barré, Gérard Brunel, Pierre Barthet, and Sara Laurencin-Dalieux. 2017. The visual analogue scale: An easy and reliable way of assessing perceived stress. *Quality in Primary Health Care* 1, 1 (2017), 1–5.
- [10] Choon Looi Bong, Kristin Fraser, and Denis Oriot. 2016. Cognitive load and stress in simulation. *Comprehensive healthcare simulation: Pediatrics* (2016), 3–17.
- [11] Guillermo Borragán, Hichem Slama, Mario Bartolomei, and Philippe Peigneux. 2017. Cognitive fatigue: A time-based resource-sharing account. *Cortex* 89 (2017), 71–84.
- [12] Minghao Cai, Genaro Rebolledo Mendez, Gisele Arevalo, Sin Sze Tang, Yalmaz Ali Abdullah, and Carrie Demmans Epp. 2024. Toward Supporting Adaptation: Exploring Affect’s Role in Cognitive Load when Using a Literacy Game. In *Proceedings of the CHI Conference on Human Factors in Computing Systems*. 1–17.
- [13] Sara Caviola, Emma Carey, Irene C Mammarella, and Denes Szucs. 2017. Stress, time pressure, strategy selection and math anxiety in mathematics: A review of the literature. *Frontiers in psychology* 8 (2017), 286709.
- [14] Eran Chajut and Daniel Algom. 2003. Selective attention improves under stress: implications for theories of social cognition. *Journal of personality and social psychology* 85, 2 (2003), 231.
- [15] Samantha WT Chan, Shardul Sapkota, Rebecca Mathews, Haimo Zhang, and Suranga Nanayakkara. 2020. Prompt: Investigating receptivity to prompts based on cognitive load from memory training conversational agent. *Proceedings of the ACM on interactive, mobile, wearable and ubiquitous technologies* 4, 4 (2020), 1–23.
- [16] Xiuli Chen, Aditya Acharya, Antti Oulasvirta, and Andrew Howes. 2021. An adaptive model of gaze-based selection. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*. 1–11.
- [17] Xiuli Chen, Gilles Bailly, Duncan P Brumby, Antti Oulasvirta, and Andrew Howes. 2015. The emergence of interactive behavior: A model of rational menu search. In *Proceedings of the 33rd annual ACM conference on human factors in computing systems*. 4217–4226.
- [18] Yi-Chun Chen, Mykel J Kochenderfer, and Matthijs TJ Spaan. 2018. Improving offline value-function approximations for POMDPs by reducing discount factors. In *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 3531–3536.
- [19] Nicolas Chopin, Omiros Papaspiliopoulos, et al. 2020. *An introduction to sequential Monte Carlo*. Vol. 4. Springer.
- [20] Andy Cockburn, Per Ola Kristensson, Jason Alexander, and Shumin Zhai. 2007. Hard lessons: Effort-inducing interfaces benefit spatial learning. In *Proceedings of the SIGCHI conference on Human factors in computing systems*. 1571–1580.
- [21] Ton De Jong. 2010. Cognitive load theory, educational research, and instructional design: Some food for thought. *Instructional science* 38, 2 (2010), 105–134.
- [22] E Greef De Tjerck, FR Arciszewski Henryk, and Mark A Neerincx. 2010. Adaptive automation based on an object-oriented task model: Implementation and evaluation in a realistic c2 environment. *Journal of Cognitive Engineering and Decision Making* 4, 2 (2010), 152–182.
- [23] Cary Deck, Salar Jahedi, and Roman Sheremeta. 2021. On the consistency of cognitive load. *European Economic Review* 134 (2021), 103695.

- [24] B Dharani and TV Geetha. 2013. Adaptive learning path generation using colored Petri nets based on behavioral aspects. In *2013 International conference on recent trends in information technology (ICRTIT)*. IEEE, 459–465.
- [25] Michael Dorneich, Stephen Whitlow, Patricia May Ververs, Jim Carciofini, and Janet Creaser. 2004. Closing the loop of an adaptive system with cognitive state. In *Proceedings of the Human Factors and Ergonomics Society Annual Meeting*, Vol. 48. SAGE Publications Sage CA: Los Angeles, CA, 590–594.
- [26] Michael C Dorneich, Patricia May Ververs, Santosh Mathan, and Stephen D Whitlow. 2005. A joint human-automation cognitive system to support rapid decision-making in hostile environments. In *2005 IEEE International Conference on Systems, Man and Cybernetics*, Vol. 3. IEEE, 2390–2395.
- [27] Wai-Tat Fu and Peter Pirolli. 2007. SNIF-ACT: A cognitive model of user navigation on the World Wide Web. *Human-Computer Interaction* 22, 4 (2007), 355–412.
- [28] Steel Crate Games. 2015. Keep Talking and Nobody Explodes. <https://keeptalkinggame.com> [Video game].
- [29] Adrian M Haith and John W Krakauer. 2018. The multiple effects of practice: skill, habit and reduced cognitive load. *Current opinion in behavioral sciences* 20 (2018), 196–201.
- [30] Hao He and Yucheng Duan. 2026. Beyond performance: A POMDP-based machine learning framework for expert cognition. *Behavior Research Methods* 58, 1 (2026), 6.
- [31] Patrick Hemmer, Max Schemmer, Michael Vössing, and Niklas Kühl. 2021. Human-AI Complementarity in Hybrid Intelligence Systems: A Structured Literature Review. *PACIS* (2021), 78.
- [32] Patrick Hemmer, Monika Westphal, Max Schemmer, Sebastian Vetter, Michael Vössing, and Gerhard Satzger. 2023. Human-AI Collaboration: The Effect of AI Delegation on Human Task Performance and Task Satisfaction. In *Proceedings of the 28th International Conference on Intelligent User Interfaces*. 453–463.
- [33] Chin-Ming Hong, Chih-Ming Chen, Mei-Hui Chang, and Shin-Chia Chen. 2007. Intelligent web-based tutoring system with personalized learning path guidance. In *Seventh IEEE International Conference on Advanced Learning Technologies (ICALT 2007)*. IEEE, 512–516.
- [34] Andrew Howes, Xiuli Chen, Aditya Acharya, and Richard L Lewis. 2018. Interaction as an emergent property of a Partially Observable Markov Decision Process. *Computational interaction design* (2018), 287–310.
- [35] M Sazzad Hussain, Rafael A Calvo, and Fang Chen. 2014. Automatic cognitive load detection from face, physiology, task performance and fusion during affective interference. *Interacting with computers* 26, 3 (2014), 256–268.
- [36] Jung Sung Kang and Seung Jun Lee. 2022. Concept of an intelligent operator support system for initial emergency responses in nuclear power plants. *Nuclear Engineering and Technology* 54, 7 (2022), 2453–2466.
- [37] Antti Kangasrääsio, Kumaripaba Athukorala, Andrew Howes, Jukka Corander, Samuel Kaski, and Antti Oulasvirta. 2017. Inferring cognitive models from data using approximate Bayesian computation. In *Proceedings of the 2017 CHI conference on human factors in computing systems*. 1295–1306.
- [38] Caitlin Kelleher and Wint Hnin. 2019. Predicting cognitive load in future code puzzles. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*. 1–12.
- [39] David E Kieras, David E Meyer, James A Ballas, and Erick J Lauber. 2000. Modern computational perspectives on executive mental processes and cognitive control: Where to from here. *Control of cognitive processes: Attention and performance XVIII* (2000), 681–712.
- [40] Mikko Lauri, David Hsu, and Joni Pajarinen. 2022. Partially observable markov decision processes in robotics: A survey. *IEEE Transactions on Robotics* 39, 1 (2022), 21–40.
- [41] David Lindlbauer, Anna Maria Feit, and Otmar Hilliges. 2019. Context-aware online adaptation of mixed reality interfaces. In *Proceedings of the 32nd annual ACM symposium on user interface software and technology*. 147–160.
- [42] George A Miller. 1956. The magical number seven, plus or minus two: Some limits on our capacity for processing information. *Psychological review* 63, 2 (1956), 81.
- [43] Vamsi Krishna Nadimpalli, Florian Hauser, Dominik Bittner, Lisa Grabinger, Susanne Staufer, and Jürgen Mottok. 2023. Systematic Literature Review for the Use of AI Based Techniques in Adaptive Learning Management Systems. In *Proceedings of the 5th European Conference on Software Engineering Education*. 83–92.
- [44] Nargess Nourbakhsh, Fang Chen, Yang Wang, and Rafael A Calvo. 2017. Detecting users' cognitive load by galvanic skin response with affective interference. *ACM Transactions on Interactive Intelligent Systems (TiiS)* 7, 3 (2017), 1–20.
- [45] Fred GWC Paas and Jeroen JG Van Merriënboer. 1994. Instructional control of cognitive load in the training of complex cognitive tasks. *Educational psychology review* 6 (1994), 351–371.
- [46] Adam Persin and Ajay Jasr. 2016. Twisting the alive particle filter. *Methodology and Computing in Applied Probability* 18 (2016), 335–358.
- [47] Lars J Planke, Yixiang Lim, Alessandro Gardi, Roberto Sabatini, Trevor Kistan, and Neta Ezer. 2020. A cyber-physical-human system for one-to-many uas operations: Cognitive load analysis. *Sensors* 20, 19 (2020), 5467.
- [48] Raphaëlle N Roy, Stephane Bonnet, Sylvie Charbonnier, and Aurélie Campagne. 2013. Mental fatigue and working memory load estimation: interaction and implications for EEG-based passive BCI. In *2013 35th annual international conference of the IEEE Engineering in Medicine and Biology Society (EMBC)*. IEEE, 6607–6610.
- [49] Sayanti Roy, Trey Smith, Brian Coltin, and Tom Williams. 2023. I need your help... or do i? maintaining situation awareness through performative autonomy. In *Proceedings of the 2023 ACM/IEEE international conference on human-robot interaction*. 122–131.
- [50] Jason R Sattizahn, Jason S Moser, and Sian L Beilock. 2016. A closer look at who "chokes under pressure". *Journal of Applied Research in Memory and Cognition* 5, 4 (2016), 470–477.
- [51] Zhangfan Shen, Linghao Zhang, Xing Xiao, Rui Li, and Ruoyu Liang. 2020. Icon familiarity affects the performance of complex cognitive tasks. *i-Perception* 11, 2 (2020), 2041669520910167.
- [52] Gaganpreet Singh, Raphaëlle N Roy, and Caroline PC Chanel. 2022. Pomdp-based adaptive interaction through physiological computing. In *HAI2022: Augmenting Human Intellect: Proceedings of the First International Conference on Hybrid Human-Artificial Intelligence*. SAGE Publications 1 Oliver's Yard, 55 City Road, London, EC1Y 1SP, 32–45.
- [53] John Sweller. 1988. Cognitive load during problem solving: Effects on learning. *Cognitive science* 12, 2 (1988), 257–285.
- [54] Tarek Taha, Jaime Valls Miró, and Gamini Dissanayake. 2008. POMDP-based long-term user intention prediction for wheelchair navigation. In *2008 IEEE International Conference on Robotics and Automation*. IEEE, 3920–3925.
- [55] Tarek Taha, Jaime Valls Miró, and Gamini Dissanayake. 2011. A POMDP framework for modelling human interaction with assistive robots. In *2011 IEEE International Conference on Robotics and Automation*. IEEE, 544–549.
- [56] Blaise Thomson, F Jurčićek, M Gašić, Simon Keizer, François Mairesse, Kai Yu, and Steve Young. 2010. Parameter learning for POMDP spoken dialogue models. In *2010 IEEE Spoken Language Technology Workshop*. IEEE, 271–276.
- [57] Stella Tomasi, David Schuff, and Ozgur Turetken. 2018. Understanding novelty: how task structure and tool familiarity moderate performance. *Behaviour & Information Technology* 37, 4 (2018), 406–418.
- [58] Manuel Valle Torre, Catharine Oertel, and Marcus Specht. 2024. The Sequence Matters in Learning-A Systematic Literature Review. In *Proceedings of the 14th Learning Analytics and Knowledge Conference*. 263–272.
- [59] Matti Vihola. 2014. Ergonomic and reliable Bayesian inference with adaptive Markov chain Monte Carlo. *Wiley statsRef: statistics reference online* (2014), 1–12.
- [60] David L Woods, Mark M Kishiyama, E William Yund, Timothy J Herron, Ben Edwards, Oren Poliva, Robert F Hink, and Bruce Reed. 2011. Improving digit span assessment of short-term verbal memory. *Journal of clinical and experimental neuropsychology* 33, 1 (2011), 101–111.