

"AI enhances our performance, I have no doubt this one will do the same": The Placebo Effect Is Robust to Negative Descriptions of AI

Agnes M. Kloft*
Aalto University
Espoo, Finland
agnes.kloft@aalto.fi

Robin Welsch*
Aalto University
Espoo, Finland
robin.welsch@aalto.fi

Thomas Kosch
HU Berlin
Berlin, Germany
thomas.kosch@hu-berlin.de

Steeven Villa
LMU Munich
Munich, Germany
villa@posthci.com

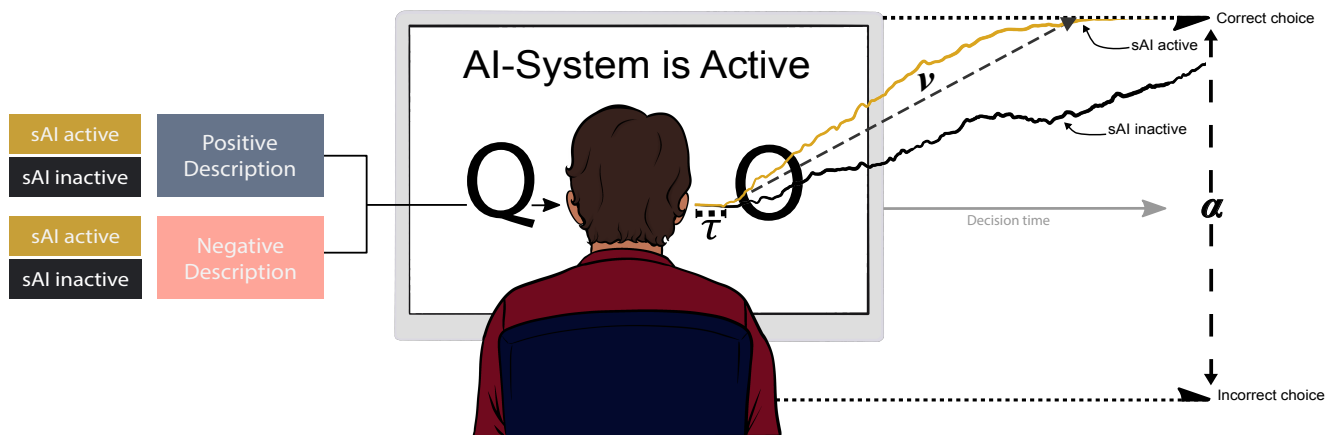


Figure 1: Schematic representation of the drift-diffusion process of decision-making with an increased drift-rate v and a decreased non-decision time τ , for the sham-AI active (sAI active) condition as compared to the sham-AI inactive (sAI inactive) condition. When using a sham AI, participants accumulate information faster.

ABSTRACT

Heightened AI expectations facilitate performance in human-AI interactions through placebo effects. While lowering expectations to control for placebo effects is advisable, overly negative expectations could induce nocebo effects. In a letter discrimination task, we informed participants that an AI would either increase or decrease their performance by adapting the interface, when in reality, no AI was present in any condition. A Bayesian analysis showed that participants had high expectations and performed descriptively better irrespective of the AI description when a sham-AI was present. Using cognitive modeling, we could trace this advantage back to participants gathering more information. A replication study verified that negative AI descriptions do not alter expectations, suggesting that performance expectations with AI are biased

*Shared first authorship: Both authors contributed equally to the paper

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).
CHI '24, May 11–16, 2024, Honolulu, HI, USA
© 2024 Copyright held by the owner/author(s).
ACM ISBN 979-8-4007-0330-0/24/05.
<https://doi.org/10.1145/3613904.3642633>

and robust to negative verbal descriptions. We discuss the impact of user expectations on AI interactions and evaluation.

CCS CONCEPTS

• **Human-centered computing** → *User studies; Empirical studies in HCI.*

KEYWORDS

Placebo, Decision-making, Performance expectation, Artificial Intelligence

ACM Reference Format:

Agnes M. Kloft, Robin Welsch, Thomas Kosch, and Steeven Villa. 2024. "AI enhances our performance, I have no doubt this one will do the same": The Placebo Effect Is Robust to Negative Descriptions of AI. In *Proceedings of the CHI Conference on Human Factors in Computing Systems (CHI '24)*, May 11–16, 2024, Honolulu, HI, USA. ACM, New York, NY, USA, 24 pages. <https://doi.org/10.1145/3613904.3642633>

1 INTRODUCTION

Expectations regarding Artificial intelligence (AI) fundamentally affect how we use this technology. The placebo effect of AI in Human-Computer Interaction (HCI) [42], inspired by medical research [4, 43, 53, 76], documents that a sham-AI (sAI) system can

bring real subjective benefits accompanied by changes in decision-making and physiology [42, 84]. Kosch et al. [42] argued that much like in the medical context, user expectations about AI technology significantly influence study outcomes and thus undermine scientific evaluation if they are left uncontrolled. The idea of controlling user expectations about novel technologies (such as AI) in human-centered studies has been discussed in the past [6], proposing to control results that originate from participant beliefs [89] rather than from an active system. Thus, user expectations play a critical role in assessing AI systems, regardless of a functional system in user studies.

Prior research on placebo effects in HCI has been reported in various contexts. For example, in gaming, where fake power-up elements that make no difference to gameplay [18] and sham descriptions of AI adaptation increase game immersion [17]. In social media, sham control settings for a news feed can result in higher user satisfaction [80]. A study by Kosch et al. [42] showed that expecting benefits from using an adaptive AI can improve subjective performance. Additionally, Villa et al. [84] could show that high expectations regarding sAI-based augmentation systems increase risky decision-making and affect information processing. Thus, AI technology can induce placebo effects that alter subjective performance, decision-making, and therefore experiences through heightened positive expectations. It is important to mention that placebo studies as the aforementioned, where a control condition is compared to a placebo, differ from placebo-controlled studies, where an effective treatment is compared to a placebo condition¹. Note, that although conceptually similar to Wizard of Oz paradigms frequently employed in AI research, where an experimenter operates a computer to simulate an intelligent system², in placebo studies, the AI system is not functional.

There are three major shortcomings in the placebo literature in HCI for AI. First, direct effects on a behavioral level are yet to be found [42, 84]. Second, it is unclear whether nocebo effects (low expectations impairing behavior) are equally influential as positive expectations based on verbal descriptions in HCI. Third, we lack scientific studies that show how AI expectations affect interaction and, thus, study outcomes.

This paper investigates the antecedents and consequences of AI's placebo effect in HCI. In detail, we examine how descriptions can impact decision-making by raising or lowering expectations, thus using expectations as a mediator between descriptions and placebo or nocebo effects. In an experimental study ($N = 66$), we examined the influence of negative and positive verbal AI descriptions and analyzed the impact of expectations on decision-making in a letter discrimination task, with and without a sAI system.

First, in line with Kosch et al. [42], Villa et al. [84], we found a subjective placebo effect: participants upheld positive expectations for the sAI system's effectiveness post-interaction. Second, we observed a main effect at the behavioral level. A Bayesian cognitive model of decision-making revealed that participants gathered information faster and altered their response style, giving us granular insights into which aspects of interaction are affected by the placebo effect. Third, contrary to previous work [18, 42, 84], we

found no effect of verbal descriptions. Participants were biased, expecting increased performance with AI, irrespective of the verbal descriptions (AI performance bias). We replicate this bias in an online study ($N = 95$).

Our results resonate with the power of AI narratives [12, 13, 39, 65] and recent calls in HCI to control for placebo effects in evaluating AI systems [42, 84]. We add an AI performance bias to the literature, which makes the AI's placebo effect robust to manipulations of verbal system descriptions. By utilizing a cognitive model of decision-making, we describe which aspects of interaction are affected by the placebo effect. We also discuss how, in a human-centered design process, the evaluation of AI must be done with user expectations in mind.

2 RELATED WORK

2.1 Expectations and the placebo effect of AI

People hold expectations with regard to AI. Survey findings show that fears about AI's disruptive impact outweigh excitement in the British public [12, 13]. This aligns with Sartori et al.'s report on the prevalence of 'AI anxiety' over perceived benefits [65]. Interestingly, it appears that the prevalence of concerns may also be influenced by narratives. For instance, science fiction portrayals have been suggested to contribute to the observed imbalance [32]. The narratives about AI can differ among stakeholders and change over time [7]. Indeed, national policies in countries like China, Germany, the USA, and France underscore AI's disruptive potential [2], and these narratives are widely impactful, affecting usage [7, 39]. Prior studies have explored key areas like transparency expectations [52, 56, 59], human-AI relationships [92], trust [21, 49, 79, 85], and autonomy [52], forming the basis for AI interface design. However, there is a gap in understanding expectations of human-AI interaction outcomes, such as task performance with AI support [42]. To address this gap, it is important to understand how exactly user expectations influence the outcome in human-AI interaction and to investigate the role narratives play in this.

The placebo effect relies on expectations [3, 19, 34, 43, 58, 63] and is not confined to medical contexts but also penetrates performance contexts like sports [5]. Here, an inert substance given to athletes can improve but also deteriorate performance [36]. While placebo effects of AI in HCI and their effect on performance have recently been studied [17, 42, 80, 84], there is very little knowledge on nocebo effects. In HCI, a nocebo effect would negatively affect both performance and subjective metrics, like usability or user experience [42]. For example, Halbhuber et al. [28] manipulated latency descriptions in gaming, showing that negative expectations reduced performance and user experience. In human-AI interaction, Ragot et al. [60] found that AI-generated art labeled as such was rated less favorably than if labeled as human-made. Thus, although first studies indicate the possibility of nocebo effects brought upon by technological artifacts, empirical studies directly leveling or even implementing negative expectations for AI are scarce. Likewise, it is unclear whether system descriptions as put forth by Kosch et al. [42] determine AI expectations which lead to placebo effects, or whether placebo effects are driven by general biases as in Ragot et al. [60]. While the former could be addressed in a study context, the latter could only be addressed within a societal discourse

¹For a taxonomy of the placebo effect in HCI see Kosch et al. [42]

²See Dahlbäck et al. [15], Schoonderwoerd et al. [67]

[12, 13]. **Therefore, it is critical to study how descriptions of AI influence placebo effects in HCI evaluation.**

2.2 Decision-making with AI

Decision-making, a process shaped by expectations and perceptions, involves selecting from a range of options [73]. The Drift Diffusion Model (DDM) serves as a cognitive framework for understanding this process, describing it as evidence accumulation until a decision boundary (a correct vs. an incorrect answer) is reached [47, 54, 61, 62]. In its most basic form, the DDM models reaction times based on correct and incorrect responses in a random walk process toward a decision boundary, see Figure 1. For a binary decision task with equal probability, we can assume three parameters. A drift-rate v , indicating the speed of gathering information, a boundary separation α , reflecting a decision-strategy, and a non-decision time τ parameter, reflecting motor preparation and perceptual processes unrelated to decision-making [47]. This model has been successful in predicting decision-making under uncertainty and in different cognitive tasks [61, 62]. Indeed, recent research argues that computational cognitive models like the DDM are central for interaction (see Oulasvirta et al. [55]). In line with this, the DDM has been applied to pedestrian crossing [91], moving target selection [45], interactions with robots [35] or teleoperations [14]. Recent studies indicate that even sham adaptive AIs can influence user performance and risk-taking in decision-making [42, 84]. However, the cognitive mechanisms behind these effects remain unclear. Applying the DDM could potentially shed light on the cognitive basis of the placebo effect for adaptive AI systems. Considering previous studies by Kosch et al. [42] and Villa et al. [84], it appears plausible that the decision criterion may be affected by the implementation of positive expectations (placebo) improving performance (more liberal decision-making with decreased α) and negative expectations (nocebo), resulting in an enlargement of α . **Consequently, users may make rapid, less accurate decisions when aided by an adaptive AI interface or slower, more accurate decisions when the AI system potentially hampers their performance.**

3 RESEARCH MODEL

We conducted a mixed-design lab study with one between- and one within-subject factor, each with two levels. Two groups (between-subject) with different system descriptions, referred to as DESCRIPTION ("the AI system worsens performance and increases stress," referred to as NEGATIVE VERBAL DESCRIPTION condition vs. "the AI system enhances performance and decreases stress," referred to as POSITIVE VERBAL DESCRIPTION condition) were investigated. The within-subject factor for each group was the sAI SYSTEM STATUS (sAI ACTIVE condition vs. sAI INACTIVE). The ORDER³ of conditions in SYSTEM STATUS was counterbalanced across participants in both DESCRIPTION groups.

³ORDER was treated as a within-subject factor in the statistical analysis (Section 6.4), addressing the question, "Is this condition from the first or the second half of the experiment?"

4 METHOD

In the following, we motivate and document our methodological choices in realizing the study. The analysis with all associated measures can be found at osf.io. The study was pre-registered, and the pre-registration details can be accessed at: aspredicted.org. Deviations from the pre-registration can be found in Table 6.

4.1 Verbal Description

The study introduction varied in its verbal DESCRIPTION between two groups. Participants in the NEGATIVE VERBAL DESCRIPTION group were informed that the AI system had "**decreased task performance**" and resulted in an "**elevation in stress**" among first users. Moreover, they were informed that the system was new and untried, thus making it "**unreliable**" and "**risky**" for use in real-world scenarios. In contrast, the participants in the POSITIVE VERBAL DESCRIPTION group were informed that the system had previously "**enhanced task performance**" while "**reducing stress**." They were also informed that the system was "**cutting-edge**," "**reliable**" and "**safe**" to use in real-world scenarios (see Appendix A for full descriptions). We informed all participants that they would be testing an AI system under two conditions: once with the AI SYSTEM STATUS set to active (sAI ACTIVE condition) and once inactive (sAI INACTIVE condition). For the sAI ACTIVE condition, participants were informed that the AI system was continuously adapting the task difficulty based on their task performance and stress levels, monitored through electrodermal activity via electrodes (see Appendix B). In contrast, in the sAI INACTIVE condition, participants were informed that the AI system was not active and that the task pace was random (see Appendix C).

4.2 Measures

4.2.1 Letter discrimination task. Two-alternative forced choice tasks, such as letter discrimination tasks, model simple decision-making and its underlying cognitive processes [48, 61, 78, 86]. In the task, participants must identify which of two letters, displayed on either side of a central target letter, matches the target. We used four letter pairs (E/F, P/R, C/G, Q/O), selected from Ratcliff and Rouder [61]. Each trial consisted of a three-component trial sequence, which began with a fixation cross centrally displayed between the letters for a variable time (interstimulus interval, ISI), facilitating perception of the system's adaptability similar to [84]. Then, one of the letters was shown for 50.1 ms in the center of the screen [78]. After this, a randomly sketched line mask rotated by $x \cdot 360$ degrees ($x \in [0,1[$) and mirrored (vertically and/or horizontally, or neither) was shown in place of the target letter for 1500 ms, see Figure 2. During the line mask presentation, the participants responded by pressing the left or the right arrow key with their index and middle finger. The first key press response during mask presentation time was recorded. The only critical change made to the task of Thapar et al. [78] was the randomly varying ISI. This was done to allow participants to track potential changes related to adaptation and should not affect task performance.

Each participant underwent 400 trials derived from two Blocks \times 100 trials of one random letter pair \times two SYSTEM STATUS conditions (sAI ACTIVE vs. sAI INACTIVE). The order of the SYSTEM STATUS conditions was counterbalanced across the participants in each

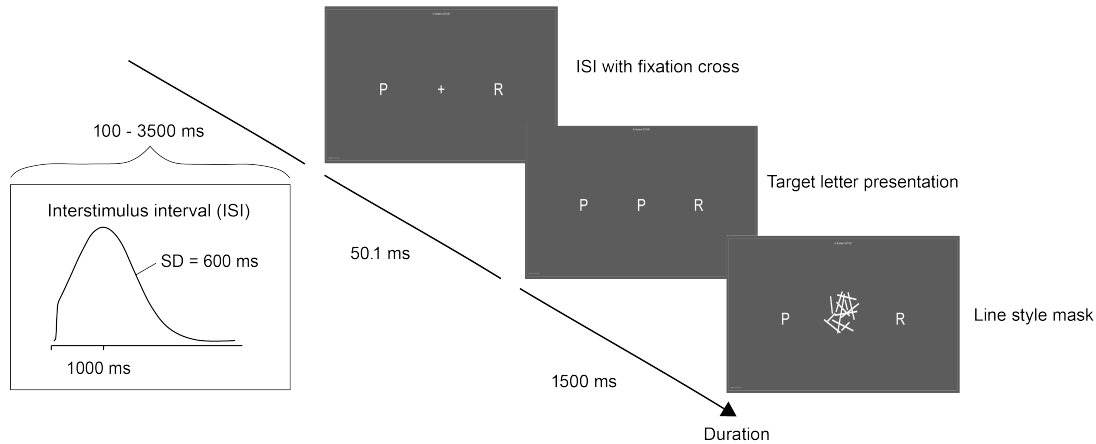


Figure 2: Trial sequence during the letter discrimination task. The duration of the ISI followed a Gaussian distribution ($M = 1000$ ms, $SD = 600$ ms). Key responses (left or right arrow) were logged during the presentation of the mask.

DESCRIPTION group. The duration of each trial varied based on the randomized duration of the ISI in the trial sequence, which followed a Gaussian distribution with a mean of 1000 ms and a standard deviation of 600 ms. The shortest trial lasted 1650.1 ms, the longest lasted 5050.1 ms, and the median trial duration was 2550.1 ms. The overall median task duration for all 400 trials was approximately 17 minutes. After each block, participants were offered a short break.

4.2.2 Questionnaires.

Assessment of expectations. We measured user expectations of performance and how they persisted after the interaction. For overall performance expectations (judgments prior to interaction), we used four questions: A seven-point Likert item (1: Strongly disagree to 7: Strongly agree), indicating the expected overall performance (*I think I will perform better in the task with the AI system than in the task without the AI system.*), a slider item from zero (slower) to 100 (faster) as an indicator for the subjectively estimated task speed (*I will be [slower/faster] in the task with the AI system active than in the task with the AI system inactive.*), and two open text questions (allowed response range: 0 to 100) asking participants the expected number of correct letter discriminations in each condition (*Out of 200 actions, how many do you expect to get correct [with/without] the AI system active?*). To evaluate judgments of performance following the interaction, identical questions phrased in the past tense were assessed. An additional questionnaire adapted from Villa et al. [84] was termed "System evaluation" and implemented to assess the participant's judgment of performance after the interaction, see Table 2.

Task load. To measure workload, we implemented the NASA-TLX [30], a well-established questionnaire [41], with six dimensions: mental demand, physical demand, temporal demand, performance, effort, and frustration. Participants rated each dimension on a scale of 1 to 20, with higher scores indicating higher task loads. We calculated the raw score by summing up the item scores (Raw-TLX, [29]).

Additional Questionnaires. We assess user experience using the UEQ-S [68] (8 item pairs; Likert scale from -3 to +3) with its two dimensions, pragmatic quality and hedonic quality. For measuring Usability, we used an adapted version of the System Usability Scale (SUS) [8], changing "system" to "AI system," adding the synonym *awkward* for *cumbersome* [22], and computed the SUS score by summing the score contributions of each item and multiplying the sum by a factor of 2.5 in line with Brooke [8].

4.2.3 Electrodermal activity recording and pre-processing. Skin conductance, reflecting physiological arousal, was measured as an indicator for cognitive workload [41] following the framework for Electrodermal Activity (EDA) research in HCI [1]. EDA was recorded using standard Ag/AgCl electrodes (24 mm surface diameter) placed on the distal surfaces of the proximal phalanges of the index and middle fingers of the participant's non-dominant hand. Before testing, participants washed their hands with soap and cleaned the areas where the electrodes were placed using a 70% alcohol wipe. For data acquisition, we used the BITalino biomedical toolkit [27] to acquire the EDA signals via Bluetooth connection. The *OpenSignals (r)evolution* Python API Version 1.2.6⁴ was set at a sampling rate of 100 Hz. Time-series data were recorded using the Lab Streaming Layer (LSL)⁵. For offline data preprocessing, we used the Neurokit toolbox [51]. After non-negative deconvolution analysis, we derived one metric of physiological arousal: the mean tonic SCL in each block.

4.3 Participants

Participants were recruited through print advertisements in the Helsinki metropolitan area. Eligibility criteria included: no background in computer science, age above 18, self-reported normal or corrected-to-normal vision, no silver allergy, and no use of medication or history of epilepsy or other cognitive/motor impairments.

⁴<https://github.com/BITalinoWorld/revolution-python-api#bitalino-revolution-python-api/>

⁵<https://github.com/labstreaminglayer/>

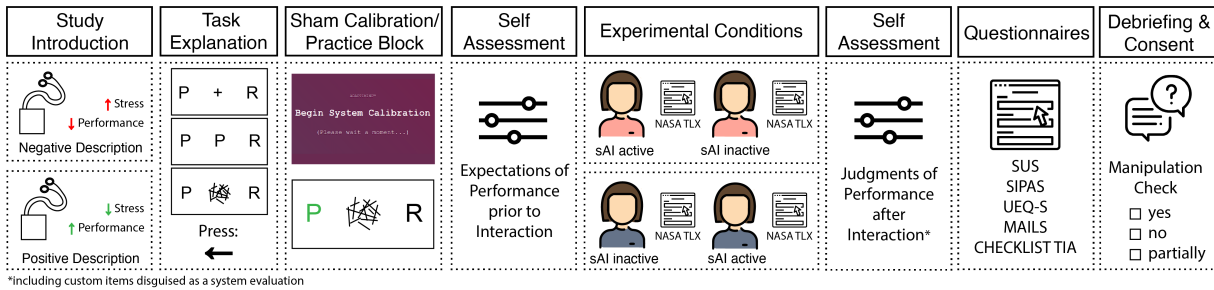


Figure 3: Study Procedure. In this mixed-design study examining the induction of placebo and nocebo effects, participants were divided into two groups (DESCRIPTION), with each group receiving altered system descriptions (NEGATIVE: AI decreased task performance and increased stress in users/ POSITIVE: AI increased task performance and decreased stress in users). Participants in each group performed a letter discrimination task under two conditions (SYSTEM STATUS): in the SHAM-AI (sAI) ACTIVE condition, they were informed that an AI system was active and adjusting the task pace based on their measured stress responses; in the sAI INACTIVE condition, they were told that the AI system was inactive and adjustments in task pace were random. The ORDER of SYSTEM STATUS alternated within each DESCRIPTION group. Before and after interacting with the sAI system, expectations on performance with and without the sAI system set as active were assessed. After the tasks and before debriefing, additional questionnaires assessing i.e., user experience and AI literacy were implemented. Ultimately, we revealed the manipulation and assessed the participants’ belief in the manipulation.

The participants received 20 Euro S-ryhmä gift vouchers as compensation for their participation. The study was approved by an ethics committee (Grant Nr. "D/594/03.04/2023").

We tested 66 participants in our study⁶, excluding one for insufficient English proficiency and one for careless responding (i.e. responding consistently with the maximum on a scale). Our final sample size consisted of 64 participants ($N = 64$, male = 24, female = 40, zero non-binary or did not disclose) with an average age of 27.64 years ($SD = 6.49$; min = 18; max = 49) reporting an average technical competence of 4.80 ($SD = 1.25$) on a 1 (low) - 7 (high) Likert item. To ensure that the two samples (DESCRIPTION: $n_{\text{positive}} = 31$, $n_{\text{negative}} = 33$) are comparable, we checked their AI literacy using the Meta AI Literacy Scale⁷ (MAILS) [10], the Checklist for Trust between People and Automation (TiA) [37] and the Subjective Information Processing Awareness Scale (SIPAS) [69–71]. We indeed found no differences as a function of DESCRIPTION, see Table 5.

4.4 Procedure

After consenting in line with the Declaration of Helsinki, the Bitlino device’s electrodes were attached, and the device was activated and secured. The experimental program appeared on the screen. We then collected data on age, profession, handedness, caffeine or medication use, experimenter familiarity, and technical competence.

Participants read an introductory text explaining the AI system and apparatus, see Figure 3. Depending on the DESCRIPTION assignment, the text included a positive or negative verbal description (Section 4.1) before interacting with the sAI. This was followed by a survey asking for information on the system being evaluated, see Villa et al. [84].

Before the task, participants completed 50 practice trials with visual feedback labeled as calibration. We then assessed their performance expectations with and without the AI system set to active. Next, participants performed the task, starting with either the sAI ACTIVE or sAI INACTIVE condition⁸, depending on the assigned ORDER. Task load was evaluated post-condition using TLX [30]. After both conditions, the AI system was further assessed (Section 4.2.2). Participants were then debriefed before re-consenting, and their belief in the manipulation was checked (Section 6.1). Finally, they were thanked and compensated for their participation. The entire experiment lasted approximately 70 minutes.

4.5 Bayesian Data Analysis and Inference

We adopted a Bayesian approach, utilizing Bayesian linear mixed models⁹. For parameter estimation, we used brms [9], a wrapper for the STAN-sampler [11] executed in R [77]. Two Hamilton-Monte-Carlo chains were computed, each with 8,000-40,000 iterations and a 20% warm-up. Trace plots of the Markov-chain Monte-Carlo permutations were inspected for divergent transitions and autocorrelation, and we checked for local convergence. All Rubin-Gelman statistics [25] were well below 1.1 and the Effective Sampling Size was over 1000. Model specifications and their non-informative priors alongside all estimated parameters can be found in Appendix H.

We then analyzed the posterior of the model. To investigate a parameter’s distinguishability from zero, we utilized p_b , which resembles the classical p -value but quantifies the effect’s likelihood of being zero or opposite [33, 74]. Effects with $p_b \leq 2.5\%$ were deemed distinguishable. We also calculated the 95% High-Density Interval (HDI) for each model parameter. For population-level effects in simple regression models, we set priors for regression parameters

⁶Deviation from pre-registration see Table 6

⁷We only implemented items of factors loading onto the dimension "AI Literacy"

⁸During the entire task, information was displayed on the screen indicating that the AI SYSTEM STATUS was set to either active or inactive

⁹For a guide on Bayesian techniques, see [9, 20, 38, 66, 82]

to one standard deviation of the outcome variable. All binary factors were effect coded (TIME (pre/post): 1, -1; SYSTEM STATUS (sAI ACTIVE/sAI INACTIVE): 1, -1; DESCRIPTION (negative/positive): 1, -1); ORDER (first condition/ second condition): 1, -1.

4.6 Apparatus

The experiment was carried out using Chromium on a Linux (Ubuntu 22.04.2 LTS) laptop (Dell Latitude 7310) with an i5 (Intel Core i5-1031U) processor and 16GB of RAM. A separate monitor (HP E27uG4) displayed the paradigm with a screen size of 27 inches (2160px*1440px) and a refresh rate of 60 Hz. The monitor's position was adjusted according to the participant's eye level. Screen distance was roughly 60 cm (23,6 inch). We built a custom experiment that ran locally with JavaScript using the lab.js library version 20.2.4 [31].

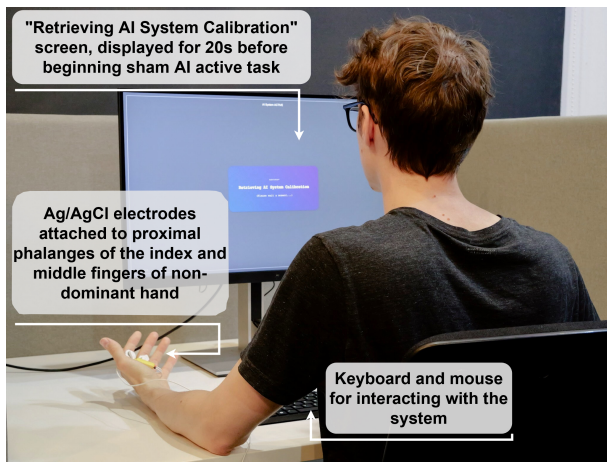


Figure 4: The participants interacted with the system with their dominant hand using a keyboard and a mouse.

5 RESEARCH QUESTIONS AND HYPOTHESES

We address the following research questions and hypotheses:

RQ1: Do subjective ratings on performance and mental workload differ between negative and positive verbal descriptions (nocebo/placebo)?

- **H1:** Lower subjective performance (**H1.1**) and higher mental workload (**H1.2**) in sAI active with a negative description (nocebo) compared to sAI inactive.
- **H2:** Higher subjective performance (**H2.1**) and lower mental workload (**H2.2**) in sAI active with a positive description (placebo) compared to sAI inactive.

RQ2: Do verbal descriptions of a sAI affect decision-making (e.g., in a letter discrimination task)?

- **H3:** More conservative speed-accuracy trade-off in sAI active with a negative description (nocebo) compared to sAI inactive.
- **H4:** More liberal speed-accuracy trade-off in sAI active with a positive description (placebo) compared to sAI inactive.

RQ3: Do verbal descriptions of a sAI affect physiological indicators of cognition when compared to no implementation of a sAI?

- **H5:** Higher levels of physiological arousal (measured by mean tonic skin conductance level (SCL)) in sAI active with a negative description (nocebo) compared to sAI inactive.
- **H6:** Lower levels of physiological arousal (measured by mean tonic SCL) in sAI active with a positive description (placebo) compared to sAI inactive.

6 RESULTS

6.1 Manipulation Check

To the question *Did you believe that an AI system was implemented to adapt task pace?* with possible answers being *Yes*, *No* or *Partially*, 10 of 64 (15.62%; 6 of 33 negative description; 4 of 31 positive description) responded with "no" and did not believe in the system's capabilities. 27 out of 64 participants (42.19%; 13 for positive description, 14 for negative description) participants reported some suspicion of the system's functionality. Thus, 27 of 64 participants fully believed in the implemented system.

6.2 Performance Expectations and Judgments of Performance

6.2.1 Subjective overall performance. To analyze expected overall performance, we centered the values by subtracting four points of the Likert item so that 0 indicates not favoring any condition and modeled overall performance estimates as a function of TIME and DESCRIPTION¹⁰. Overall, participants were positive about the sAI, $Intercept = 0.51 [0.25, 0.77]$, $p_b = 0.00\%$. However, participants showed no difference in subjective performance before and after interaction with the sAI ($\tilde{b}_{Time} = 0.19 [-0.03, 0.42]$, $p_b = 6.70\%$). There was no main effect of DESCRIPTION ($\tilde{b}_{Description} = -0.23 [-0.50, 0.03]$, $p_b = 4.66\%$) and no interaction effects ($\tilde{b}_{Time \times Description} = 0.16 [-0.06, 0.38]$, $p_b = 7.99\%$), see also Figure 5A.

6.2.2 Subjective estimated task speed. We computed a similar model to investigate the participants' expected task speed by subtracting 50 points so that zero indicates a neutral response. Figure 5B shows the average expected speed across all conditions being positive, $Intercept = 8.54 [5.41, 11.78]$, $p_b = 0.00\%$. The participants believed to be faster with the sAI active before interacting with the system ($M = 62.47$, $SD = 17.31$) than after ($M = 54.56$, $SD = 18.01$). This difference ($d_z = 0.30$) could be distinguished from zero, $\tilde{b}_{Time} = 3.96 [0.92, 7.03]$, $p_b = 0.50\%$. We found no differences for DESCRIPTION ($\tilde{b}_{Description} = -1.31 [-4.52, 1.88]$, $p_b = 21.09\%$) or interaction effects, $DESCRIPTION \times TIME$ $\tilde{b}_{Description \times Time} = -1.16 [-4.17, 1.93]$, $p_b = 22.46\%$.

6.2.3 Subjective estimated number of correct responses. We expanded the statistical model to consider SYSTEM STATUS for estimated points (no transformation) in each condition; see also Figure 5C. Participants indicated that in the sAI ACTIVE condition ($M = 142.46$, $SD = 35.87$), they would score more points than in the sAI INACTIVE condition ($M = 129.77$, $SD = 36.56$). This difference was not zero $\tilde{b}_{System Status} = 6.33 [3.23, 9.26]$, $p_b = 0.00\%$, $d_z = 0.53$.

¹⁰Gaussian link-function with default priors.

Participants believed to score more points before performing the task ($M = 142.36$, $SD = 36.16$) than after ($M = 129.88$, $SD = 33.31$, $d_z = 0.47$), $\tilde{b}_{\text{Time}} = 6.31$ [3.31, 9.33], $p_b = 0.00\%$, resembling Kosch et al. [42]. We found no distinguishable effects for DESCRIPTION $\tilde{b}_{\text{Description}} = -1.35$ [-8.78, 5.78], $p_b = 35.51\%$, or any interaction effects $p_b > 4.07\%$, see also Figure 5C.

Therefore, participants were biased toward a superior performance with AI even when given a negative verbal description of the system. We refer to this as AI PERFORMANCE BIAS.

6.3 Performance data

We excluded 6 out of 64 participants (9.38%) only from the behavioral data analysis as they did not comply with our task (percent correct <60% in one of the conditions or very large number of misses >35%). We deleted the first trial in each trial block along with too-short responses by filtering reaction times (RT) under 150 ms (519 out of 23084; 2.25%)¹¹ and missed responses with RT > 1499 ms (32 out of 22565; 0.14%).

To explore our interventions, we computed two separate regression models with varying intercepts for each participant and, ORDER (first vs second experimental block), SYSTEM STATUS and DESCRIPTION as population-level effects for RT (Gaussian-link function) and correctness of response (Bernoulli-link function). For RT, we found an effect for SYSTEM STATUS, $\tilde{b}_{\text{System Status}} = -4.17$ ms [-6.14, -2.17], $p_b = 0.00\%$. Participants reacted on average faster in the sAI ACTIVE condition ($M = 604$ ms, $SD = 92$ ms) to stimuli as compared to the sAI INACTIVE condition ($M = 611$ ms, $SD = 79$ ms; Cohen's $d_z = 0.12$) Figure 6A. We also found that participants increased their response speed from the first to the second experimental condition; we found an effect for ORDER, $\tilde{b}_{\text{Order}} = 11.05$ ms [9.06, 13.06], $p_b = 0.00\%$ (First condition: $M = 619$ ms, $SD = 91$ ms; Second condition: $M = 596$ ms, $SD = 79$ ms; Cohen's $d_z = 0.39$). There was no effect of DESCRIPTION, $\tilde{b}_{\text{Description}} = -4.76$ ms [-26.43, 16.43], $p_b = 32.80\%$. For the correctness of responses, we found the same pattern of results. Participants were more likely to respond correctly in the sAI ACTIVE ($M = 90.07\%$, $SD = 9.20\%$) condition as compared to the sAI INACTIVE condition ($M = 89.35\%$, $SD = 8.80\%$; $\tilde{b}_{\text{System Status}} = -0.05$ [0.00, 0.09], $p_b = 2.05\%$; Odds = 0.95) and improved in accuracy along the course if the experiment (ORDER), $\tilde{b}_{\text{Order}} = -0.05$ [-0.09, 0.00], $p_b = 1.37\%$, Odds = 1.05 (First condition: $M = 89.29\%$, $SD = 9.75\%$; Second condition: $M = 90.13\%$, $SD = 8.18\%$). There was no effect of DESCRIPTION, $\tilde{b}_{\text{Description}} = 0.10$ [-0.12, 0.33], $p_b = 19.00\%$. For descriptives of the performance data as a function of SYSTEM STATUS and DESCRIPTION, see Table 1.

We computed the DDM¹² to test H3 & H4 on the reaction time data¹³, see Figure 6A. A hierarchical form of this model was built accounting for inter-subject variability with a varying intercept and a population-level effect for each SYSTEM STATUS and an interaction term for DESCRIPTION for each τ , ν and α .

¹¹Deviation from pre-registration see Table 6

¹²Keep in mind that in the DDM we model RTs based on correct and incorrect responses by fitting data to a model that represents decision-making as the noisy accumulation of information (ν denoting the average rate of accumulation), for one choice or the other, until a threshold is reached (α ; boundary separation). A starting point from which the accumulation process starts and a parameter τ denoting non-decision time is added to the model. For a visual representation, see Figure 1

¹³Deviation from pre-registration see Table 6

We inspected the parameters of the model, see Figure 6B, for differences in SYSTEM STATUS Figure 6B-10 and DESCRIPTION Figure 6B-11 for boundary separation α . See Figure 7A, to see whether the difference in reaction time and percent correct comes from a change in the participant's strategy, e.g., prioritizing speed over the accuracy. We found that in the sAI active condition, participants had a slightly larger boundary separation, α , making them slightly more conservative as compared to the sAI inactive condition (Figures 6B-10 and 7A). However, we also found that ν (drift rate), see Figure 7B, was higher for sAI active as compared to sAI inactive. Thus, information accumulation was relatively faster in the sAI active condition, see Figure 6B-6. With a relatively faster accumulation of information, ν , and more conservative boundaries, α , in the sAI active condition as compared to sAI inactive, we can explain the differences between conditions for the singular analysis of RT and the correctness of trials (for a schematic representation of this difference for SYSTEM STATUS, see Figure 1). Note that when we visually inspected the posterior distribution for each participant, as well as their RT difference as a function of SYSTEM STATUS (Appendix F), we found that the effect did not vary as a function of post-experimental belief in the system, see Section 6.1. Therefore, the model seems to hold for all participants, irrespective of their beliefs after debriefing.

Similarly, τ was also affected by SYSTEM STATUS with an interaction with DESCRIPTION qualifying the effect, see Table 12. Looking at Figure 6A and Figure 7C, we can see that the group with the negative description had a slightly earlier onset in RT. For all parameter values, see Table 12 and for the mathematical formulation and priors, see Appendix E. To contextualize the effect size on RT, we also predicted the RT from the model, averaged across conditions, and calculated Cohen's d_z for ORDER, at 1.21, and for SYSTEM STATUS, at 0.74.

6.4 Workload and physiological arousal

Investigating H1.2, H2.2, H5 and H6, we computed a regression model for the NASA-TLX raw data with SYSTEM STATUS, DESCRIPTION, their interaction and ORDER as predictors found. We found no differences for SYSTEM STATUS $\tilde{b}_{\text{System Status}} = -0.08$ [-2.04, 1.90], $p_b = 46.82\%$, DESCRIPTION $\tilde{b}_{\text{Description}} = 0.83$ [3.72, 5.42], $p_b = 35.91\%$, their interaction effects $\tilde{b}_{\text{System Status} \times \text{Description}} = 0.43$ [-1.54, 2.39], $p_b = 33.28\%$ or ORDER, $\tilde{b}_{\text{Order}} = 1.45$ [-0.52, 3.42], $p_b = 7.28\%$. For EDA,¹⁴ there was no effect of the SYSTEM STATUS, $\tilde{b}_{\text{System Status}} = -0.20$ [-0.53, 0.13], $p_b = 11.88\%$, no effect of DESCRIPTION $\tilde{b}_{\text{Description}} = 0.25$ [-0.10, 0.59], $p_b = 7.37\%$, no interaction effect, $\tilde{b}_{\text{Description} \times \text{System Status}} = 0.07$ [-0.24, 0.38], $p_b = 32.70\%$ or ORDER effect, $\tilde{b}_{\text{Order}} = 0.20$ [-0.03, 0.42], $p_b = 4.20\%$.

6.5 Usability and User Experience

Except for *The AI system made the task easier* (item 2), which was viewed more favorably with a positive description, there were no significant differences in DESCRIPTION (Table 2). Participants slightly disagreed with *The task was easy* (item 1) and were slightly negative about *The AI system improved my cognitive abilities* (item

¹⁴Same predictor formula; 6 participants excluded due to poor signal quality

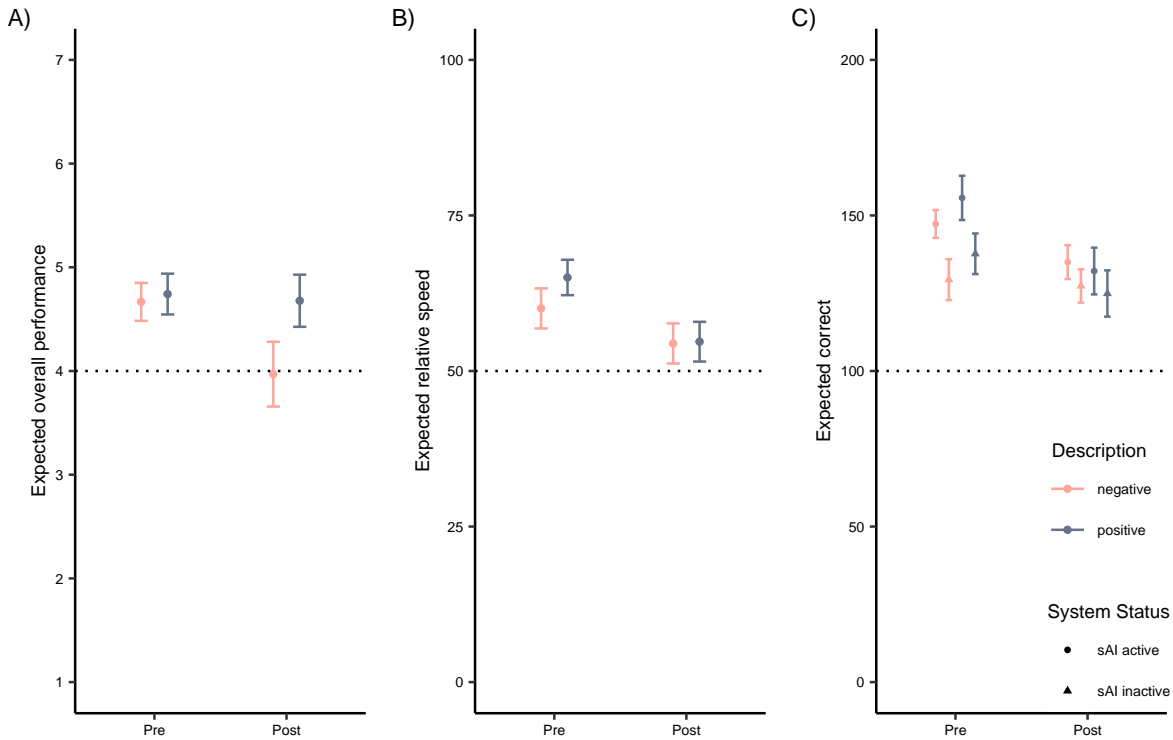


Figure 5: A: Mean expected performance as a function of TIME and DESCRIPTION. B: Mean expected relative speed as a function of TIME and DESCRIPTION. C: Mean expected correct responses before and after interacting with the sAI system as a function of TIME, SYSTEM STATUS and DESCRIPTION. Error bars denote ± 1 standard error of the mean.

Table 1: Mean percent correct and reaction time (RT) for both correct and incorrect trials as a function of SYSTEM STATUS and DESCRIPTION

Description	sAI active			sAI inactive		
	Correct %	Correct RT	Incorrect RT	Correct %	Correct RT	Incorrect RT
Negative	91.11 (8.86)	586.80 (75.61)	719.16 (168.01)	90.35 (6.68)	596.56 (43.71)	739.32 (156.84)
Positive	88.87 (9.61)	599.44 (96.61)	723.33 (147.19)	88.21 (10.76)	603.19 (95.69)	722.31 (161.11)

7). Yet, similar to Kosch et al. [42], Villa et al. [84], they agreed that the AI has future potential.

For UEQ-S scales, we found an overall positive user experience, with no group differences on hedonic or pragmatic attributes, with both having positive values indicating a positive user experience. SUS ratings indicated that the system was rated average in terms of usability unaffected by DESCRIPTION.

7 REPLICATION STUDY: POSITIVE EXPECTATIONS FOR NEGATIVE DESCRIPTIONS

To confirm the AI PERFORMANCE BIAS, we conducted an additional online replication study with the negative system description. We

replicated the first part of the previous study using the same negative verbal description and subjective questions to assess expectations and judgments. Subsequently, we replaced the adaptation description, which initially referred to utilizing real-time EDA analysis to measure stress responses, with the use of computer vision technology to analyze facial expressions in real-time, as per Kosch et al. [42] (no data was recorded). To address potential concerns about participants not fully comprehending the instructions, we set up the experiment to enforce comprehension of verbal descriptions. Based on this, the participants were divided into two groups. Both groups read the negative system description. However, one group was asked to complete a comprehension check (COMPREHENSION), ensuring they fully understood the negative description, before

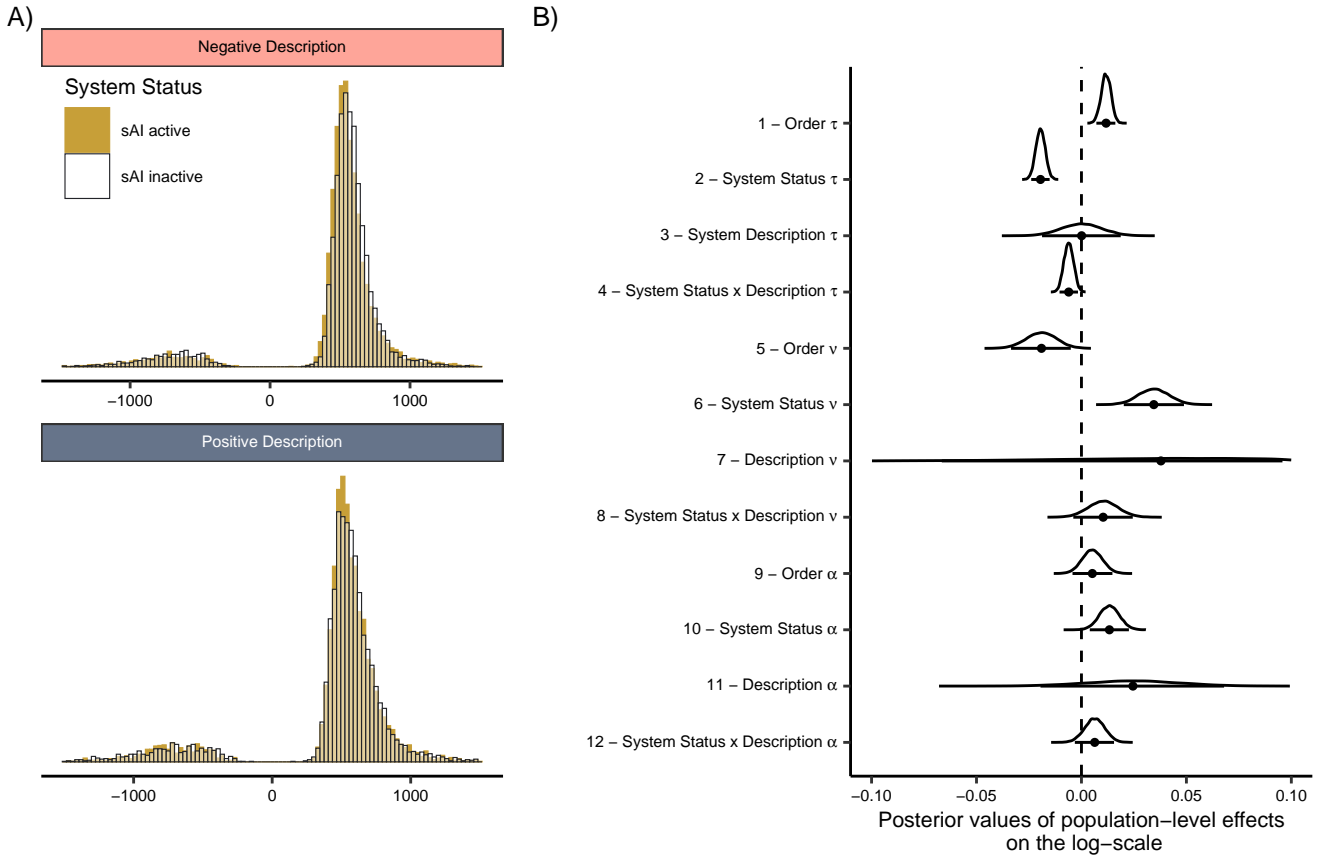


Figure 6: A: Reaction time distribution as a function of SYSTEM STATUS (sAI active vs. sAI inactive) and DESCRIPTION (incorrect trials are multiplied with -1). B: Posterior density plot for the parameter values for all population-level parameters 95% High-Density Interval (HDI). If the HDI does not cross the midline p_b will be <2.5%.

being able to continue to the next part of the study. In the NO-COMPREHENSION group, participants were not bound by the same requirement, allowing for variations in their engagement with the negative description. This decision was made to facilitate a comparison between participants in the comprehension group, where individuals were required to fully understand the text, indicating a predicted decline in performance, and the no-comprehension group. While some may not have read it thoroughly, others may have held pre-existing expectations. This contrast allows a nuanced exploration of how differing levels of understanding might influence participants’ responses.

We recruited 95 participants via prolific. Five participants had to be excluded due to incomplete data, e.g., missing responses in demographics, or too short or incomprehensible responses to open questions, leaving 90 participants (Age: $M = 30.69$, $SD = 9.17$, $Min = 18$, $Max = 65$) for analysis. The first group ($N_{No-Comprehension} = 44$) completed the check and got no feedback on correctness, while the second group ($N_{Comprehension} = 46$) had to answer all questions correctly (coded no: -1/yes: 1) to continue with the study. After the check the participants gave their assessment of how they expected to perform with the AI system. Finally, participants explained their point choices in an open text field. The study took, on average,

about 10 minutes to complete. Participants were compensated at £13.48/hr, resulting in a payment of £2.25 for a 10 minute-survey.

7.1 Quantitative results

Table 3 shows all means of the subjective performance expectations for each group. COMPREHENSION had an effect on overall performance, $\tilde{b}_{Comprehension} = -0.26 [-0.48, -0.04]$, $p_b = 1.05\%$ and expected task speed, $\tilde{b}_{Comprehension} = -4.60 [-8.34, -0.86]$, $p_b = 0.82\%$. For estimated correct, we added SYSTEM STATUS to the model. COMPREHENSION had no effect $\tilde{b}_{Comprehension} = -0.51 [-7.80, 6.74]$, $p_b = 44.50\%$. However, a difference for SYSTEM STATUS emerged, $\tilde{b}_{System\ Status} = 5.71 [2.00, 9.39]$, $p_b = 0.15\%$ and an interaction effect $\tilde{b}_{System\ Status \times Comprehension} = -4.36 [-8.08, -0.69]$, $p_b = 1.03\%$. Participants in the group without the enforced comprehension check estimated to answer more accurately with the sAI active than without ($p_b\ diff = 0.00\%$), while in the comprehension check group, this difference was not present ($p_b\ diff = 30.40\%$). Most importantly, participants were optimistic with regard to overall performance and expected speed, irrespective of COMPREHENSION. Only for the difference in the number of expected correct responses, we find that the COMPREHENSION leveled participants to neutral expectations.

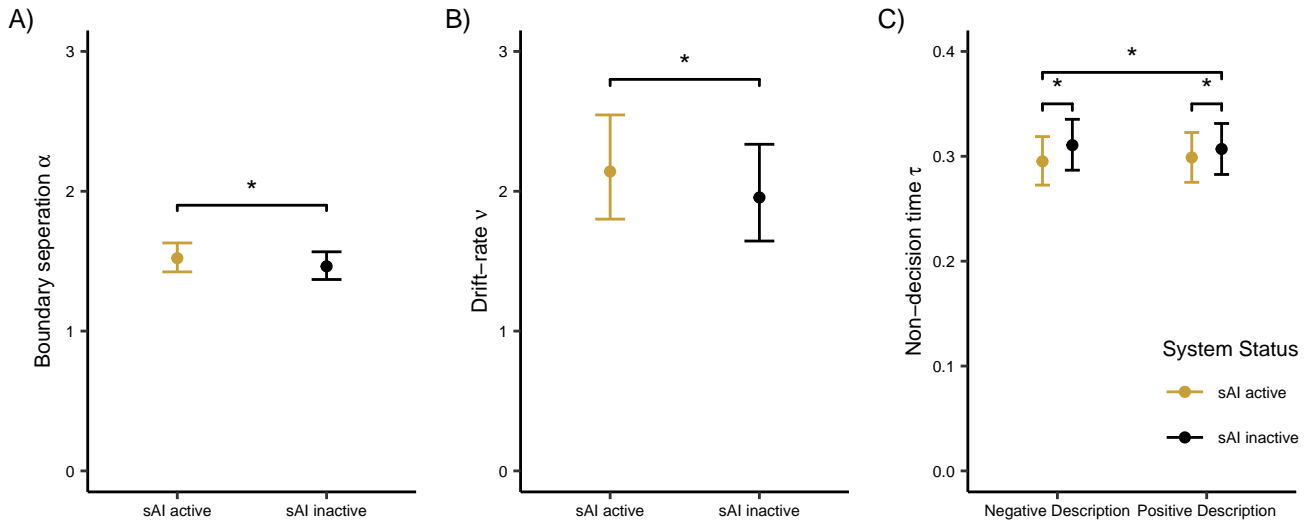


Figure 7: A: Estimates with High-Density Interval (HDI) 95% for boundary separation, α , as a function of SYSTEM STATUS. B: Estimates with HDI 95% for the drift rate ν as a function of SYSTEM STATUS. C: Estimates of non-decision time τ with HDI 95% as a function of SYSTEM STATUS and DESCRIPTION.

Table 2: The customized system evaluation was answered on nine 7-point Likert items (1: strongly disagree; 7: strongly agree). We estimate the difference towards a neutral value and compare the samples across DESCRIPTION. A neutral value for the custom items was 4; for the UEQ-S scales, it was set to zero. We fitted a robust regression model for each comparison. For the adapted SUS, the expected average is 68. Distinguishable effects from a neutral value (expected) or for each DESCRIPTION are marked with *. We used a studentized link-function with priors scaled to one SD

Item/scale	M_{neg} (SD)	M_{pos} (SD)	$\tilde{\Delta}_{expected}$ [HDI 95%]	p_b	$\tilde{b}_{Description}$ [HDI 95%]	p_b
System evaluation						
1 - The task was easy.	3.24 (1.62)	3.61 (1.82)	-0.60 [-1.03,-0.16]	0.45%*	-0.17 [-0.63, 0.25]	21.58%
The AI system						
2 - made the task easier.	3.36 (1.78)	4.19 (1.42)	-0.24 [-0.66, 0.16]	12.54%	-0.41 [-0.82, 0.00]	2.48%*
3 - made the task more enjoyable.	3.33 (1.83)	4.03 (1.43)	-0.32 [-0.74, 0.10]	6.39%	-0.35 [-0.76, 0.07]	5.15%
4 - made me more self-confident.	3.39 (1.49)	3.97 (1.72)	-0.33 [-0.74, 0.08]	5.87%	-0.27 [-0.68, 0.14]	9.35%
5 - made me more efficient.	3.87 (1.55)	4.16 (1.27)	0.02 [-0.33, 0.39]	24.80%	-0.12 [-0.48, 0.23]	44.47%
6 - improved my performance.	3.96 (1.55)	4.19 (1.33)	0.09 [-0.27, 0.45]	31.37%	-0.11 [-0.46, 0.27]	28.33%
7 - improved my cognitive abilities.	3.55 (1.39)	3.77 (1.34)	-0.35 [-0.68, -0.01]	2.06%*	-0.09 [-0.43, 0.25]	29.47%
8 - has a lot of potential for future development.	5.03 (1.15)	5.42 (1.12)	1.23 [0.93, 1.51]	0.00%*	-0.20 [-0.49, 0.09]	8.98%
UEQ-S-Pragmatic	0.54 (0.93)	0.85 (1.19)	0.73 [0.46, 0.99]	0.00%*	-0.17 [-0.43, 0.10]	10.79%
UEQ-S-Hedonic	0.73 (1.03)	0.80 (1.32)	0.78 [0.49, 1.08]	0.00%*	-0.04 [-0.33, 0.26]	40.80%
SUS-Adapted	64.62 (13.86)	67.74 (17.43)	-1.41 [-5.39, 2.47]	23.55%	-1.71 [-5.69, 2.18]	19.28%

Note: In this figure, "neg" denotes a negative system description, while "pos" represents a positive one.

7.2 Qualitative results

After participants estimated their subjective performance, they were further prompted to elaborate on the rationale behind their responses. To gain deeper insights into participants' perceptions

and expectations regarding their performance with AI, a qualitative analysis was performed. The focus was on revising statements made by participants regarding their expectations of performing better or worse with an AI system when informed of a potential performance

Table 3: Summary statistics for performance expectations as a function of COMPREHENSION

Performance expectations with sAI active when compared to sAI inactive	Comprehension	
	M_{No} (SD)	M_{Yes} (SD)
Overall performance*	5.11* (1.13)	4.59* (0.98)
Task speed*	70.05* (18.48)	60.70* (17.19)
Difference in number of correct responses	-20.16* (29.71)	-2.72 (39.41)

Note: Differences between groups are highlighted in the variable with a *. Means that are distinguishably more positive than their neutral value (4 for overall performance, 50 for task speed and zero for Difference in the number for correct responses) are marked with *.

decline (NEGATIVE DESCRIPTION). This qualitative exploration aimed to uncover nuanced reasons underlying the participants' convictions about performance with AI and the perceived speed advantage or disadvantage.

The analysis involved clustering statements based on the participants' subjective assessments of their expected speed and overall performance on the Likert items. Two researchers independently performed a qualitative analysis of the statements, grouping them according to their semantic meaning. Afterward, a consensus was reached, identifying five distinct categories: AI Trust, AI Assistance, Uncertainty, Neutral, Self-Awareness, and AI Antagonism.

Table 4 provides a detailed breakdown of the distribution of statements across these categories, revealing predominant themes. Notably, the majority of statements (out of 180) primarily align with AI Trust (27 statements), AI Assistance (64 statements), and Uncertainty (44 statements). AI Trust reflects the participants' positive expectations and trust in the capabilities of AI systems as powerful tools that ensure an advantage. AI Assistance describes the perception of AI as a helpful assistant that facilitates task completion. Uncertainty portrays the participants' uncertainty toward the AI system's influence on task completion. These prevalent themes indicate that the majority of participants expected a positive influence (AI Trust and AI Assistance) on task performance ($N_{Statements} = 41$) and speed ($N_{Statements} = 50$) from the AI system, with some expressing uncertainty instead of negative sentiment toward the AI system despite being informed of potential performance decline.

8 DISCUSSION

In this study, we set out to implement negative expectations and study the nocebo effect of AI (RQ1). However, we found that the placebo effect of AI in HCI [42] is robust to the manipulation of expectations by a negative verbal description (contrary to H1.1 and H1.2). Even when we told participants that the AI system would make the task harder and more stressful, they still believed it would improve their performance. This was the same for those who read positive descriptions of the AI (rendering H1, H3 & H5 void). We refer to this expectation of high performance as AI PERFORMANCE BIAS. We replicate this bias in a dedicated online study.

We found that heightened expectations (supporting H2.1.) carry over to the way participants make decisions (RQ2). Participants in the sAI active condition responded slightly faster and more accurately when informed they were interacting with an adaptive AI system. Using the DDM model to analyze decision-making, we found that believing an AI is involved can make participants gather information more quickly, respond more conservatively, and make

them more alert (partial support for H4). We found no effects on workload or physiological arousal (no support for H2.1, H6).

8.1 Beyond demand characteristics and system descriptions

Critics may argue that placebo effects in AI are not genuine and stem from demand characteristics, which often influence experimental studies and HCI evaluations [16, 88]. In our study, despite participants being primed to view the AI negatively, their improved performance and positive ratings contradicted these expectations, suggesting that demand characteristics cannot account for the AI's placebo effect. One could also assert that our system descriptions were not effective in producing expectations. We used similar positive and negative verbal descriptions as studies in sports science, e.g., [5, 36]. Also, the manipulation of SYSTEM STATUS influenced participants both subjectively and behaviorally, irrespective of their post-experimental accounts of believing in the system's capabilities, see Section 6.1. Moreover, in Study 2, participants who understood the negative AI description (comprehension check) adjusted their expectations accordingly. This indicates that while our negative portrayal had some impact, it was less influential than AI narratives, that created high expectations. Future research should further explore this by comparing a sham AI system with a non-AI system (e.g., controlled by a sham operator) or by screening for AI expectations a-priori and comparing the placebo response with a rather neutral and minimal AI description.

8.2 AI Performance Bias as an Antecedent of the Placebo Effect of AI

It appears that the prevailing positive perceptions about AI are influential enough to overshadow context-specific negative verbal descriptions, irrespective of reported belief after the experiment. This could be due to participants bringing their daily experiences and narratives of AI into the evaluation, biasing both their subjective evaluations and behavior, see Table 4. From a mental model perspective [90], performance-reducing AI assistance may not fit into a coherent representation of human-AI interaction. It follows that the placebo mechanism for AI interfaces presented in the HCI literature is invalid [42, 84], as they focus on verbal system descriptions producing a placebo effect of AI. Based on our qualitative data, we follow that the effect is not specific to verbal descriptions of the system but may arise out of the socio-technical context as a function of the user's mental model.

Table 4: Subjective influence of the AI system on expected performance and speed before task completion: Number of statements and percentage per category

Category	Description	Statement Examples	Statement Counts (%)	
			Performance	Speed
AI Trust	Trust and positive expectations toward the capabilities of AI systems in general. Seeing AI as a powerful tool that ensures an advantage.	<i>AI models enhance our performances, so I have no doubt that this one will do the same.</i> (22P; P = 6) <i>Because I trust AI, IT IT FAST [sic] and quite reliable for most activities.</i> (11S; S = 68)	15 (16.67%)	12 (13.33%)
AI Assistance	AI is a helpful assistant that will facilitate task completion.	<i>I think the AI will assist me as it will be programmed to do the task, and I am not.</i> (35P; P = 6); <i>With the help of AI, I will be able to work fast because it will be assisting me rather than having to figure this out myself.</i> (37S; S = 85); <i>My effort and AI combined we will produce better results.</i> (40S; S = 99)	26 (28.89%)	38 (42.22%)
Uncertainty	Uncertainty toward the AI's systems influence on task completion.	<i>I don't know what to expect, really, maybe I could do better or not.</i> (18P; P = 4); <i>[...] I do not know the effects of the AI system on my performance yet.</i> (56S; S = 51)	29 (32.22%)	15 (16.67%)
Neutral	AI will neither have a positive or negative influence. AI won't make a difference in the task.	<i>I don't think there will be a large effect either way.</i> (2P; P = 4); <i>Because AI shouldn't have an effect on how I respond.</i> (15S; S = 56)	7 (7.78%)	9 (10.00%)
Self-Awareness	Self-reliance, and confidence in individual abilities, regardless of AI assistance, emphasizing autonomy and individual skill.	<i>Because I do not depend on enhancement to complete my tasks.</i> (7P; P = 6); <i>Cause I am a bit smarter for now than the AI system.</i> (39S; S = 99),	9 (10.00%)	7 (7.78%)
AI Antagonism	Lack of trust in the AI system, believing it will hamper performance, and skepticism towards AI's usefulness.	<i>As far as I understand, the AI will confuse me more than be of any help.</i> (60P; P = 4); <i>The AI might distract me and make me a little slower.</i> (9S; S = 27)	4 (4.44%)	9 (10.00%)

Note: The participants' statements on the AI systems influence on their performance and speed were grammatically corrected to ensure good readability. Any quotes that remain unchanged are marked with [sic]. Each quote is followed by parentheses indicating the statement item number and whether the statement is related to the participants' assessments of their expected performance (P) or speed (S). The number after the semicolon indicates the participants' subjective assessments of their expected performance on a Likert item ranging from 1 (strongly disagree) to 7 (strongly agree). Similarly, for expected speed, participants provided scores on a scale ranging from 1 (slower) to 100 (faster).

The AI performance bias presents an intriguing contrast with Sartori and Bocca [65] findings on AI Anxiety. While individuals often express strong negative attitudes about AI replacing them in certain tasks, it appears that when humans and AI work together, even in a non-functional AI setting, joint performance is judged to be superior. Past studies have demonstrated that task performance in human-AI collaborations can surpass individual AI or human performance [23]. However, our findings shed new light on these findings. The human-AI performance gain may not arise from the summation of individual capabilities but also involves an elevation in human performance influenced by performance expectations.

This suggests that (HCI) designers may harness the advantages of human-AI collaboration when focusing on systems that leverage a symbiotic relationship rather than fully automated tasks. However, future studies should explore not only the context of collaboration similar to Villa et al. [84] and Kosch et al. [42] but also consider human-AI competition.

8.3 The Impact of Sham-AI on Decision-making

Villa et al. [84] explored the impact of the placebo effect on decision-making in risky situations. They found that individuals with high expectations of AI system support tended to take greater risks

compared to those without AI assistance. This emphasizes how people’s actions can be shaped by the narrative surrounding AI systems. In our study, we extended this research by investigating how positive and negative verbal descriptions affect decision-making processes. Our model showed that when people believed to have AI support, they gathered information faster than when not supported by AI. Yet, the type of narrative (positive or negative) did not have an impact on parameters in the DDM and, thus, the underlying decision-making process. Prior research indicates that a participant’s confidence can substantially influence the drift rate in a DDM [46, 50]. Therefore, it is possible that our findings can be explained by the participants feeling more confident when using the AI system. Also, we find a slightly more conservative decision boundary, with participants gathering more information until making a decision when supported with sAI. With AI support, participants might prioritize accuracy (a strategy that can be experimentally induced [75]), which also improves their overall performance. Lastly, sAI also shortened participants’ non-decision time, indicating they were in a more prepared state when making decisions, especially for negative descriptions. Note, however, that while some proponents associate a reduced non-decision time with better attention, as argued by Nunez et al. [54], or disinhibition [72], others have developed models without this parameter [83], as it is sensitive to contaminants. Thus, our computational model shows that the belief in using AI influenced participants’ decision-making processes when interacting with a computing system.

8.4 Limitations & Implications

The study presents multiple limitations. While fostering a comfortable and friendly environment is commonly recommended in HCI evaluations [44, 64], prior research [24] has indicated that positive emotions can counteract the nocebo response in pain experiments. Positive affect could explain why we observed no nocebo effects. Analysis of EDA and TLX data over time showed that participants, at the very least, were not strained by the task. Nonetheless, future research should take into account the impact of emotions during tests, perhaps by deliberately altering them, as suggested by Geers et al. [24].

It is worth noting that in addition to positive affect possibly accounting for the absence of nocebo effects, the fact that only around 17% of participants didn’t fully believe in the AI system’s capabilities could also serve as an explanation. However, this percentage is lower than the number of participants who either fully believed in or had some level of suspicion towards the system. Yet, the effect was present in most nonbelievers nevertheless (see Appendix F).

In line with van Berkel and Hornbæk [81], we highlight two major domains of implications of our work. First, methodologically, given that a drift rate in the DDM can be estimated fast [86], the DDM could be used to compute a robust behavioral indicator of a placebo response for an AI interface. Second, it is crucial for the HCI community to understand that technology narratives can significantly bias AI performance expectations to the point where even negative descriptions cannot mitigate their influence on evaluation and interaction. For instance, positive expectations (placebo) may lead to overconfidence regarding the attributes of

the system, such as its usability or user experience [42]. Our findings demonstrate that individuals tended to be overly confident about their performance. This could potentially mislead those evaluating the technology, fundamentally undermining the principles of human-centered design. One could argue that our behavioral effects are small and, thus, the placebo effect of AI is negligible to human-centered design. We will outline why these small effects are unproblematic regarding our claims. First, while our behavioral effects were small ($d_z = 0.12$), and arguably they become larger when controlling for the speed-accuracy trade-off and thus accounting for individual variation, effects on subjective measures were medium-sized ($d_z = 0.53$). Second, we used minimal intervention by only describing a sham AI system. A more severe intervention, including more placebo characteristics (for an overview, see [58]) may yield more substantial effects. In the context of a user study, a false-positive due to placebo could have severe consequences (for a discussion, see Kosch et al. [42]). Prentice and Miller [57] argue that small effects in studies with minimal interventions are particularly meaningful, much like Götz et al. [26] that posit how small effects are essential to progress in science. Third, placebo/nocebo interventions in sports contexts are also tied to small effects ([36] $d_{placebo} = 0.36$, $d_{nocebo} = 0.37$). Note also that studies on aging populations with similar tasks only find medium effects [62]. Given the medium-sized subjective effects that align with our small behavioral results, we deem our results meaningful for applied contexts.

8.5 Potential Strategies to Mitigate the Placebo Effect of AI Technologies

Building on previous studies demonstrating a placebo effect in HCI [42, 84], our research investigated the impact of positive or negative descriptions of AI in eliciting a placebo or nocebo effect. Contrary to our hypotheses, we were unable to induce a nocebo effect (negative descriptions leading to the expectation of a poorer performance) with AI technology. Even when AI is framed negatively, people expect it to be effective and improve performance. Based on these findings, we propose strategies for mitigating the potential influence of prior expectations when evaluating AI technologies, which should be investigated in future research:

- (1) **Monitor Decision-Making Processes:** Observe changes in participants’ judgments or behaviors in response to negative/positive information about the system, utilizing subjective, behavioral, and psychophysiological measures [40, 84].
- (2) **Minimize AI Disclosure:** Refrain from informing users about the AI’s involvement to avoid biasing their experiences and thus control for contextual placebo-related information [87].
- (3) **Transparent AI Disclosure When Necessary:** If AI disclosure is unavoidable, clearly communicate its limitations and development status to encourage critical evaluation based on performance rather than expectations.
- (4) **Incorporate Sham Conditions:** Use a non-functional AI (sham) condition alongside the functional AI in experiments to differentiate the AI’s actual effect from user expectations.
- (5) **Evaluate Expectation Narratives:** Conduct interviews to understand user anticipations and perceptions regarding

specific technologies to see how pre-existing expectations influence the study outcome.

9 CONCLUSION

We found that even when we told participants to expect poor performance from a fake AI system, they still performed better and responded faster, showing a robust placebo effect. Contrary to previous work, this indicates that the placebo effect of AI is not easily negated by negative verbal descriptions, which raises questions about current methods for controlling for expectations in HCI studies. Additionally, the belief in having AI assistance facilitated decision-making processes, even when the narrative about AI was negative, thereby emphasizing that the influence of AI goes beyond simple narratives. This highlights the complexity and impact of AI narratives and suggests the need for a more nuanced approach in both research and practical user evaluation of AI.

ACKNOWLEDGMENTS

We thank Otso Haavisto for coding the experiment, testing participants, designing recruitment posters, piloting, and commenting on the draft. We also thank Salla Nicholls for testing participants, preprocessing data, creating the Qualtrics surveys, and recruiting and scheduling participants. Last but not least, our thanks go to Jasper Quinn and Beatriz Mello for their support in piloting and participant recruitment.

REFERENCES

- [1] Ebrahim Babaei, Benjamin Tag, Tilman Dingler, and Eduardo Velloso. 2021. A Critique of Electrodermal Activity Practices at CHI. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems* (Yokohama, Japan) (CHI '21). Association for Computing Machinery, New York, NY, USA, Article 177, 14 pages. <https://doi.org/10.1145/3411764.3445370>
- [2] Jascha Bareis and Christian Katzenbach. 2022. Talking AI into Being: The Narratives and Imaginaries of National AI Strategies and Their Performative Politics. *Science, Technology, & Human Values* 47, 5 (May 2022), 855–881. <https://doi.org/10.1177/01622439211030007>
- [3] Ernest Edward Beckham. 1989. Improvement after evaluation in psychotherapy of depression: Evidence of a placebo effect? *Journal of Clinical Psychology* 45, 6 (Nov. 1989), 945–950. [https://doi.org/10.1002/1097-4679\(198911\)45:6<945::aid-jclp2270450620>3.0.co;2-2](https://doi.org/10.1002/1097-4679(198911)45:6<945::aid-jclp2270450620>3.0.co;2-2)
- [4] Henry K. Beecher. 1955. The Powerful Placebo. *Journal of the American Medical Association* 159, 17 (Dec. 1955), 1602–1606. <https://doi.org/10.1001/jama.1955.02960340022006>
- [5] Christopher J Beedie, Damian A Coleman, and Abigail J Foad. 2007. Positive and Negative Placebo Effects Resulting from the Deceptive Administration of an Ergogenic Aid. *International Journal of Sport Nutrition and Exercise Metabolism* 17, 3 (June 2007), 259–269. <https://doi.org/10.1123/ijsnem.17.3.259>
- [6] Walter R. Boot, Daniel J. Simons, Cary Stothart, and Cassie Stutts. 2013. The Pervasive Problem With Placebos in Psychology: Why Active Control Groups Are Not Sufficient to Rule Out Placebo Effects. *Perspectives on Psychological Science* 8, 4 (2013), 445–454. <https://doi.org/10.1177/1745691613491271>
- [7] Paolo Bory. 2019. Deep new: The shifting narratives of artificial intelligence from Deep Blue to AlphaGo. *Convergence* 25, 4 (Feb. 2019), 627–642. <https://doi.org/10.1177/1354856519829679>
- [8] John Brooke. 1996. SUS: A 'Quick and Dirty' Usability Scale. In *Usability Evaluation in Industry* (1st ed.), Patrick W. Jordan, B. Thomas, Ian Lyall McClelland, and Bernard Weerdmeester (Eds.). CRC Press, London, Chapter 21, 189–194. <https://doi.org/10.1201/9781498710411>
- [9] Paul-Christian Bürkner. 2017. brms: An R Package for Bayesian Multilevel Models Using Stan. *Journal of Statistical Software* 80, 1 (Aug. 2017), 1–28. <https://doi.org/10.18637/jss.v080.i01>
- [10] Astrid Carolus, Martin Koch, Samantha Straka, Marc Latoschik, and Carolin Wienrich. 2023. MAALS – Meta AI Literacy Scale: Development and Testing of an AI Literacy Questionnaire Based on Well-Founded Competency Models and Psychological Change- and Meta-Competencies. *Computers in Human Behavior: Artificial Humans* 1, 2 (Aug. 2023), 10 pages. <https://doi.org/10.1016/j.chbah.2023.100014>
- [11] Bob Carpenter, Andrew Gelman, Matthew D. Hoffman, Daniel Lee, Ben Goodrich, Michael Betancourt, Marcus Brubaker, Jiqiang Guo, Peter Li, and Allen Riddell. 2017. Stan: A Probabilistic Programming Language. *Journal of Statistical Software* 76, 1 (Jan. 2017), 1–32. <https://doi.org/10.18637/jss.v076.i01>
- [12] Stephen Cave, Kate Coughlan, and Kanta Dihal. 2019. "Scary Robots": Examining Public Responses to AI. In *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society* (Honolulu, HI, USA) (AI/ES '19). Association for Computing Machinery, New York, NY, USA, 331–337. <https://doi.org/10.1145/3306618.3314232>
- [13] Stephen Cave and Kanta Dihal. 2019. Hopes and fears for intelligent machines in fiction and reality. *Nature Machine Intelligence* 1 (Feb. 2019), 74–78. <https://doi.org/10.1038/s42256-019-0020-9>
- [14] Javier Corredor, Jorge Sofrony, and Angelika Peer. 2017. Decision-Making Model for Adaptive Impedance Control of Teleoperation Systems. *Institute of Electrical and Electronics Engineers (IEEE) Transactions on Haptics* 10, 1 (Jan. 2017), 5–16. <https://doi.org/10.1109/toh.2016.2581807>
- [15] Nils Dahlbäck, Arne Jönsson, and Lars Ahrenberg. 1993. Wizard of Oz studies: why and how. In *Proceedings of the 1st International Conference on Intelligent User Interfaces* (Orlando, FL, USA) (IUI '93). Association for Computing Machinery, New York, NY, USA, 193–200. <https://doi.org/10.1145/169891.169968>
- [16] Nicola Dell, Vidya Vaidyanathan, Indrani Medhi, Edward Cutrell, and William Thies. 2012. "Yours is better!": participant response bias in HCI. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems* (Austin, Texas, USA) (CHI '12). Association for Computing Machinery, New York, NY, USA, 1321–1330. <https://doi.org/10.1145/2207676.2208589>
- [17] Alena Denisova and Paul Cairns. 2015. The Placebo Effect in Digital Games: Phantom Perception of Adaptive Artificial Intelligence. In *Proceedings of the 2015 Annual Symposium on Computer-Human Interaction in Play* (London, United Kingdom) (CHI Play '15). Association for Computing Machinery, New York, NY, USA, 23–33. <https://doi.org/10.1145/2793107.2793109>
- [18] Alena Denisova and Elliott Cook. 2019. Power-Ups in Digital Games: The Rewarding Effect of Phantom Game Elements on Player Experience. In *Proceedings of the Annual Symposium on Computer-Human Interaction in Play* (Barcelona, Spain) (CHI Play '19). Association for Computing Machinery, New York, NY, USA, 161–168. <https://doi.org/10.1145/3311350.3347173>
- [19] Nico J. Diederich and Christopher G. Goetz. 2008. The placebo treatments in neurosciences. *Neurology* 71, 9 (Aug. 2008), 677–684. <https://doi.org/10.1212/01.wnl.0000324635.49971.3d>
- [20] Alan Dix. 2022. Bayesian statistics. In *Bayesian Methods for Interaction and Design*, John H. Williamson, Antti Oulasvirta, Per Ola Kristensson, and Nikola Banovic (Eds.). Cambridge University Press, 81–114. <https://doi.org/10.1017/9781108874830.004>
- [21] Andrea Ferrario, Michele Loi, and Eleonora Viganò. 2020. In AI We Trust Incrementally: a Multi-layer Model of Trust to Analyze Human-Artificial Intelligence Interactions. *Philosophy and Technology* 33 (Sept. 2020), 523–539. <https://doi.org/10.1007/s13347-019-00378-3>
- [22] Kraig Finstad. 2006. The System Usability Scale and Non-Native English Speakers. *Journal of Usability Studies* 1, 4 (Aug. 2006), 185–188. <https://dl.acm.org/doi/10.5555/2835531.2835535>
- [23] Andreas Fügener, Jörn Grahl, Alok Gupta, and Wolfgang Ketter. 2022. Cognitive Challenges in Human-Artificial Intelligence Collaboration: Investigating the Path Toward Productive Delegation. *Information Systems Research* 33, 2 (June 2022), 678–696. <https://doi.org/10.1287/isre.2021.1079>
- [24] Andrew L Geers, Shane Close, Fawn C Caplandies, Charles L Vogel, Ashley B Murray, Yopina Pertiwi, Ian M Handley, and Lene Vase. 2019. Testing a positive-affect induction to reduce verbally induced nocebo hyperalgesia in an experimental pain paradigm. *Pain* 160, 10 (Oct. 2019), 2290–2297. <https://doi.org/10.1097/j.pain.0000000000001618>
- [25] Andrew Gelman and Donald B Rubin. 1992. Inference from Iterative Simulation Using Multiple Sequences. *Statistical science* 7, 4 (1992), 457–472. <https://doi.org/10.1214/ss/1177011136>
- [26] Friedrich M Götz, Samuel D Gosling, and Peter J Rentfrow. 2022. Small effects: The Indispensable Foundation for a Cumulative Psychological Science. *Perspectives on Psychological Science* 17, 1 (Jan. 2022), 205–215. <https://doi.org/10.1177/1745691620984483>
- [27] José Guerreiro, Raúl Martins, Hugo Silva, André Lourenço, and Ana Fred. 2013. BITalino - A Multimodal Platform for Physiological Computing. In *Proceedings of the 10th International Conference on Informatics in Control, Automation and Robotics - Volume 2: ICINCO* (Reykjavik, Iceland), Vol. 1. SciTePress, 500–506. <https://doi.org/10.5220/0004594105000506>
- [28] David Halhuber, Maximilian Schlenzcek, Johanna Bogon, and Niels Henze. 2022. Better Be Quiet about It! The Effects of Phantom Latency on Experienced First-Person Shooter Players. In *Proceedings of the 21st International Conference on Mobile and Ubiquitous Multimedia* (Lisbon, Portugal) (MUM '22). Association for Computing Machinery, New York, NY, USA, 172–181. <https://doi.org/10.1145/3568444.3568448>
- [29] Sandra G. Hart. 2006. Nasa-Task Load Index (NASA-TLX); 20 Years Later. *Proceedings of the Human Factors and Ergonomics Society Annual Meeting* 50, 9 (Oct. 2006), 904–908. <https://doi.org/10.1177/154193120605000909>

- [30] Sandra G. Hart and Lowell E. Staveland. 1988. Development of NASA-TLX (Task Load Index): Results of Empirical and Theoretical Research. In *Human Mental Workload*, Peter A. Hancock and Najmedin Meshkati (Eds.), *Advances in Psychology*, Vol. 52. North-Holland, 139–183. [https://doi.org/10.1016/S0166-4115\(08\)62386-9](https://doi.org/10.1016/S0166-4115(08)62386-9)
- [31] Felix Henninger, Yury Shevchenko, Ulf Mertens, Pascal J. Kieslich, and Benjamin E. Hilbig. 2021. *lab.js: A free, open, online experiment builder*. <https://doi.org/10.5281/zenodo.5233003>
- [32] Isabella Hermann. 2020. Beware of fictional AI narratives. *Nature Machine Intelligence* 2, 11 (Nov. 2020), 654–654. <https://doi.org/10.1038/s42256-020-00256-0>
- [33] Herbert Hoijtink and Rens van de Schoot. 2018. Testing small variance priors using prior-posterior predictive p values. *Psychological Methods* 23, 3 (Sept. 2018), 561–569. <https://doi.org/10.1037/met0000131>
- [34] Asbjørn Hróbjartsson and Peter Christian Gøtzsche. 2001. Is the Placebo Powerless? An Analysis of Clinical Trials Comparing Placebo with no Treatment. *The New England Journal of Medicine* 344, 21 (May 2001), 1594–1602. <https://doi.org/10.1056/nejm200105243442106>
- [35] Jie Huang, Wenhua Wu, Zhenyi Zhang, and Yutao Chen. 2020. A Human Decision-Making Behavior Model for Human-Robot Interaction in Multi-Robot Systems. *Institute of Electrical and Electronics Engineers (IEEE) Access* 8 (Nov. 2020), 197853–197862. <https://doi.org/10.1109/access.2020.3035348>
- [36] Philip Hurst, Lieke Schipof-Godart, Attila Szabo, John Raglin, Florentina Hettinga, Bart Roelands, Andrew Lane, Abby Foad, Damian Coleman, and Chris Beedie. 2020. The Placebo and Nocebo effect on sports performance: A systematic review. *European Journal of Sport Science* 20, 3 (Aug. 2020), 279–292. <https://doi.org/10.1080/17461391.2019.1655098>
- [37] Jiun-Yin Jian, Ann M. Bisantz, and Colin G. Drury. 2000. Foundations for an Empirically Determined Scale of Trust in Automated Systems. *International Journal of Cognitive Ergonomics* 4, 1 (March 2000), 53–71. https://doi.org/10.1207/S15327566IJCE0401_04
- [38] Matthew Kay, Gregory L. Nelson, and Eric B. Heckler. 2016. Researcher-Centered Design of Statistics: Why Bayesian Statistics Better Fit the Culture and Incentives of HCI. In *Proceedings of the 2016 CHI Conference on Human Factors in Computing Systems* (San Jose, CA, USA) (CHI '16). Association for Computing Machinery, New York, NY, USA, 4521–4532. <https://doi.org/10.1145/2858036.2858465>
- [39] William R. King and Jun He. 2006. A meta-analysis of the technology acceptance model. *Information & Management* 43, 6 (Sept. 2006), 740–755. <https://doi.org/10.1016/j.im.2006.05.003>
- [40] Thomas Kosch, Jakob Karolus, Havy Ha, and Albrecht Schmidt. 2019. Your skin resists: exploring electrodermal activity as workload indicator during manual assembly. In *Proceedings of the ACM SIGCHI Symposium on Engineering Interactive Computing Systems* (Valencia, Spain) (EICS '19). Association for Computing Machinery, New York, NY, USA, Article 8, 5 pages. <https://doi.org/10.1145/3319499.3328230>
- [41] Thomas Kosch, Jakob Karolus, Johannes Zagermann, Harald Reiterer, Albrecht Schmidt, and Paweł W. Woźniak. 2023. A Survey on Measuring Cognitive Workload in Human-Computer Interaction. *Comput. Surveys* 55, 13s, Article 283 (July 2023), 39 pages. <https://doi.org/10.1145/3582272>
- [42] Thomas Kosch, Robin Welsch, Lewis Chuang, and Albrecht Schmidt. 2023. The Placebo Effect of Artificial Intelligence in Human-Computer Interaction. *ACM Transactions on Computer-Human Interaction* 29, 6, Article 56 (Jan. 2023), 32 pages. <https://doi.org/10.1145/3529225>
- [43] Louis Lasagna, Frederick Mosteller, John M. von Felsinger, and Henry K. Beecher. 1954. A study of the placebo response. *The American Journal of Medicine* 16, 6 (June 1954), 770–779. [https://doi.org/10.1016/0002-9343\(54\)90441-6](https://doi.org/10.1016/0002-9343(54)90441-6)
- [44] Jonathan Lazar, Julio Abascal, Simone Barbosa, Jeremy Barksdale, Batya Friedman, Jens Grossklags, Jan Gulliksen, Jeff Johnson, Tom McEwan, Loïc Martinez-Normand, Wibke Michalk, Janice Tsai, Gerrit van der Veer, Hans von Axelson, Ake Walldius, Gill Whitney, Marco Winckler, Volker Wulf, Elizabeth F. Churchill, Lorrie Cranor, Janet Davis, Alan Hedge, Harry Hochheiser, Juan Pablo Hourcade, Clayton Lewis, Lisa Nathan, Fabio Paterno, Blake Reid, Whitney Quesenbery, Ted Selker, and Brian Wentz. 2016. Human-Computer Interaction and International Public Policymaking: A Framework for Understanding and Taking Future Actions. *Foundations and Trends in Human-Computer Interaction* 9, 2 (May 2016), 69–149. <https://doi.org/10.1561/11000000062>
- [45] Byungjoo Lee, Sunjun Kim, Antti Oulasvirta, Jong-In Lee, and Eunji Park. 2018. Moving Target Selection: A Cue Integration Model. In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems* (Montreal, QC, Canada) (CHI '18). Association for Computing Machinery, New York, NY, USA, 1–12. <https://doi.org/10.1145/3173574.3173804>
- [46] Douglas G. Lee, Jean Daunizeau, and Giovanni Pezzulo. 2021. Evidence or Confidence: What Is Really Monitored during a Decision? *Psychonomic Bulletin & Review* 30 (March 2021), 1360–1379. <https://doi.org/10.3758/s13423-023-02255-9>
- [47] Veronika Lerche and Andreas Voss. 2018. Speed-accuracy manipulations and diffusion modeling: Lack of discriminant validity of the manipulation or of the parameter estimates? *Behavior Research Methods* 50 (March 2018), 2568–2585. <https://doi.org/10.3758/s13428-018-1034-7>
- [48] Veronika Lerche, Andreas Voss, and Markus Nagler. 2017. How many trials are required for parameter estimation in diffusion modeling? A comparison of different optimization criteria. *Behavior Research Methods* 49 (April 2017), 513–537. <https://doi.org/10.3758/s13428-016-0740-2>
- [49] Bingjie Liu. 2021. In AI We Trust? Effects of Agency Locus and Transparency on Uncertainty Reduction in Human-AI Interaction. *Journal of Computer-Mediated Communication* 26, 6 (Nov. 2021), 384–402. <https://doi.org/10.1093/jcmc/zmab013>
- [50] Yaxin Liu and Stella F. Lourenco. 2022. Drift diffusion modeling informs how affective factors affect visuospatial decision making. *Journal of Vision* 22, 14, Article 3394 (Dec. 2022). <https://doi.org/10.1167/jov.22.14.3394>
- [51] Dominique Makowski, Tam Pham, Zen J. Lau, Jan C. Brammer, François Lespinasse, Hung Pham, Christopher Schölzl, and S. H. Annabel Chen. 2021. NeuroKit2: A Python toolbox for neurophysiological signal processing. *Behavior Research Methods* 53 (Feb. 2021), 1689–1696. <https://doi.org/10.3758/s13428-020-01516-y>
- [52] Christian Meurisch, Cristina A. Mihale-Wilson, Adrian Hawlitschek, Florian Giger, Florian Müller, Oliver Hinz, and Max Mühlhäuser. 2020. Exploring User Expectations of Proactive AI Systems. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies* 4, 4, Article 146 (Dec. 2020), 22 pages. <https://doi.org/10.1145/3432193>
- [53] Guy Montgomery and Irving Kirsch. 1996. Mechanisms of Placebo Pain Reduction: An Empirical Investigation. *Psychological Science* 7, 3 (May 1996), 174–176. <https://doi.org/10.1111/j.1467-9280.1996.tb00352.x>
- [54] Michael D. Nunez, Joachim Vandekerckhove, and Ramesh Srinivasan. 2017. How attention influences perceptual decision making: Single-trial EEG correlates of drift-diffusion model parameters. *Journal of Mathematical Psychology* 76 (Feb. 2017), 117–130. <https://doi.org/10.1016/j.jmp.2016.03.003>
- [55] Antti Oulasvirta, Jussi P. P. Jokinen, and Andrew Howes. 2022. Computational Rationality as a Theory of Interaction. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems* (New Orleans, LA, USA) (CHI '22). Association for Computing Machinery, New York, NY, USA, Article 359, 14 pages. <https://doi.org/10.1145/3491102.3517739>
- [56] Hyanghee Park, Daehwan Ahn, Kartik Hosanagar, and Joonhwan Lee. 2021. Human-AI Interaction in Human Resource Management: Understanding Why Employees Resist Algorithmic Evaluation at Workplaces and How to Mitigate Burdens. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems* (Yokohama, Japan) (CHI '21). Association for Computing Machinery, New York, NY, USA, Article 154, 15 pages. <https://doi.org/10.1145/3411764.3445304>
- [57] Deborah A Prentice and Dale T Miller. 1992. When small effects are impressive. In *Methodological issues and strategies in clinical research* (4th ed.), Alan e. Kazdin (Ed.). American Psychological Association, 99–105. <https://doi.org/10.1037/14805-006>
- [58] Donald D. Price, Damien G. Finnis, and Fabrizio Benedetti. 2008. A Comprehensive Review of the Placebo Effect: Recent Advances and Current Thought. *Annual Review of Psychology* 59 (Jan. 2008), 565–590. <https://doi.org/10.1146/annurev.psych.59.113006.095941>
- [59] Zoe A. Purcell, Mengchen Dong, Anne-Marie Nussberger, Nils Köbis, and Maurice Jakesch. 2023. Fears about AI-mediated communication are grounded in different expectations for one's own versus others' use. [arXiv:2305.01670 \[cs.HC\]](https://arxiv.org/abs/2305.01670)
- [60] Martin Ragot, Nicolas Martin, and Salomé Cojean. 2020. AI-Generated vs. Human Artworks: A Perception Bias Towards Artificial Intelligence?. In *Extended Abstracts of the 2020 CHI Conference on Human Factors in Computing Systems* (Honolulu, HI, USA) (CHI EA '20). Association for Computing Machinery, New York, NY, USA, 1–10. <https://doi.org/10.1145/3334480.3382892>
- [61] Roger Ratcliff and Jeffrey N. Rouder. 2000. A diffusion model account of masking in two-choice letter identification. *Journal of experimental psychology. Human perception and performance* 26, 1 (Feb. 2000), 127–40. <https://doi.org/10.1037/0096-1523.26.1.127>
- [62] Roger Ratcliff and Philip L. Smith. 2010. Perceptual discrimination in static and dynamic noise: The temporal relation between perceptual encoding and decision making. *Journal of Experimental Psychology: General* 139, 1 (Feb. 2010), 70–94. <https://doi.org/10.1037/a0018128>
- [63] K. Rickels, P. T. Hesbacher, C. C. Weise, B. Gray, and H. S. Feldman. 1970. Pills and improvement: A study of placebo response in psychoneurotic outpatients. *Psychopharmacologia* 16 (Jan. 1970), 318–328. <https://doi.org/10.1007/bf00404738>
- [64] Jeffrey Rubin and Dana Chisnell. 2008. *Handbook of usability testing: How to plan, design, and conduct effective tests*. John Wiley & Sons.
- [65] Laura Sartori and Giulia Bocca. 2023. Minding the gap(s): public perceptions of AI and socio-technical imaginaries. *AI & Society* 38, 2 (March 2023), 443–458. <https://doi.org/10.1007/s00146-022-01422-1>
- [66] Daniel J Schad, Michael Betancourt, and Shravan Vasishth. 2021. Toward a principled Bayesian workflow in cognitive science. *Psychological methods* 26, 1 (Feb. 2021), 103–126. <https://doi.org/10.1037/met0000275>
- [67] Tjeerd A.J. Schoonderwoerd, Emma M. van Zoelen, Karel van den Bosch, and Mark A. Neerincx. 2022. Design patterns for human-AI co-learning: A wizard-of-Oz evaluation in an urban-search-and-rescue task. *International Journal of Human-Computer Studies* 164, Article 102831 (Aug. 2022). <https://doi.org/10.1016/j.ijhcs.2022.102831>

- [68] Martin Schrepp, Andreas Hinderks, and Jörg Thomaschewski. 2017. Design and Evaluation of a Short Version of the User Experience Questionnaire (UEQ-S). *International Journal of Interactive Multimedia and Artificial Intelligence* 4 (Dec. 2017), 103–108. <https://doi.org/10.9781/ijimai.2017.09.001>
- [69] Tim Schrills and Thomas Franke. 2021. Subjective Information Processing Awareness Scale (SIPAS). (July 2021).
- [70] Tim Schrills and Thomas Franke. 2023. How Do Users Experience Traceability of AI Systems? Examining Subjective Information Processing Awareness in Automated Insulin Delivery (AID) Systems. *ACM Trans. Interact. Intell. Syst.* (March 2023). <https://doi.org/10.1145/3588594>
- [71] Tim Schrills, Mourad Zoubir, Mona Bickel, Susanne Kargl, and Thomas Franke. 2021. Are Users in the Loop? Development of the Subjective Information Processing Awareness Scale to Assess XAI. In *Proceedings of the ACM CHI Workshop on Operationalizing Human-Centered Perspectives in Explainable AI* (Yokohama, Japan) (Chi '21). Association for Computing Machinery, New York, NY, USA.
- [72] Stefanie Schuch. 2016. Task inhibition and response inhibition in older vs. younger adults: A diffusion model analysis. *Frontiers in Psychology* 7, Article 1722 (Nov. 2016). <https://doi.org/10.3389/fpsyg.2016.01722>
- [73] Eldar Shafir, Itamar Simonson, and Amos Tversky. 1993. Reason-based choice. *Cognition* 49, 1 (1993), 11–36. [https://doi.org/10.1016/0010-0277\(93\)90034-S](https://doi.org/10.1016/0010-0277(93)90034-S)
- [74] Haolun Shi and Guosheng Yin. 2020. Reconnecting p-value and Posterior Probability under One- and Two-sided Tests. *The American Statistician* 75 (Feb. 2020), 265–275. Issue 3. <https://doi.org/10.1080/00031305.2020.1717621>
- [75] Jeffrey J Starns and Roger Ratcliff. 2010. The effects of aging on the speed-accuracy compromise: Boundary optimality in the diffusion model. *Psychology and Aging* 25, 2 (June 2010), 377–390. <https://doi.org/10.1037/a0018022>
- [76] Steve Stewart-Williams and John Podd. 2004. The placebo effect: dissolving the expectancy versus conditioning debate. *Psychological Bulletin* 130, 2 (March 2004), 324–340. <https://doi.org/10.1037/0033-2909.130.2.324>
- [77] R Core Team. 2021. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria. <https://www.R-project.org/>
- [78] Anjali Thapar, Roger Ratcliff, and Gail McKoon. 2003. A diffusion model analysis of the effects of aging on letter discrimination. *Psychology and Aging* 18, 3 (Sept. 2003), 415–429. <https://doi.org/10.1037/0882-7974.18.3.415>
- [79] Takane Ueno, Yuto Sawa, Yeongdae Kim, Jacqueline Urakami, Hiroki Oura, and Katie Seaborn. 2022. Trust in Human-AI Interaction: Scoping Out Models, Measures, and Methods. In *Extended Abstracts of the 2022 CHI Conference on Human Factors in Computing Systems* (New Orleans, LA, USA) (CHI EA '22). Association for Computing Machinery, New York, NY, USA, Article 254, 7 pages. <https://doi.org/10.1145/3491101.3519772>
- [80] Kristen Vaccaro, Dylan Huang, Motahhare Eslami, Christian Sandvig, Kevin Hamilton, and Karrie Karahalios. 2018. The Illusion of Control: Placebo Effects of Control Settings. In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems* (Montreal, QC, Canada) (CHI '18). Association for Computing Machinery, New York, NY, USA, Article 16, 13 pages. <https://doi.org/10.1145/3173574.3173590>
- [81] Niels van Berkel and Kasper Hornbæk. 2023. Implications of Human-Computer Interaction Research. *Interactions* 30, 4 (June 2023), 50–55. <https://doi.org/10.1145/3600103>
- [82] Rens van de Schoot, Sarah Depaoli, Ruth King, Bianca Kramer, Kaspar Mårtens, Mahlet G. Tadesse, Marina Vannucci, Andrew Gelman, Duco Veen, Joukje Willemssen, and Christopher Yau. 2021. Bayesian statistics and modelling. *Nature Reviews Methods Primers* 1, 1 (Jan. 2021), 1–26. <https://doi.org/10.1038/s43586-020-00001-2>
- [83] Don van den Bergh, Francis Tuerlinckx, and Stijn Verdonck. 2020. DstarM: an R package for analyzing two-choice reaction time data with the D*M method. *Behavior Research Methods* 52 (April 2020), 521–543. <https://doi.org/10.3758/s13428-019-01249-7>
- [84] Steeven Villa, Thomas Kosch, Felix Grelka, Albrecht Schmidt, and Robin Welsch. 2023. The placebo effect of human augmentation: Anticipating cognitive augmentation increases risk-taking behavior. *Computers in Human Behavior* 146, Article 107787 (Sept. 2023). <https://doi.org/10.1016/j.chb.2023.107787>
- [85] Kailas Vodrahalli, Roxana Daneshjou, Tobias Gerstenberg, and James Zou. 2022. Do Humans Trust Advice More If It Comes from AI? An Analysis of Human-AI Interactions. In *Proceedings of the 2022 AAAI/ACM Conference on AI, Ethics, and Society* (Oxford, United Kingdom) (AI/ES '22). Association for Computing Machinery, New York, NY, USA, 763–777. <https://doi.org/10.1145/3514094.3534150>
- [86] Andreas Voss, Jochen Voss, and Veronika Lerche. 2015. Assessing cognitive processes with diffusion model analyses: a tutorial based on fast-dm-30. *Frontiers in Psychology* 6 (March 2015). <https://doi.org/10.3389/fpsyg.2015.00336>
- [87] Tor D Wager and Lauren Y Atlas. 2015. The neuroscience of placebo effects: connecting context, learning and health. *Nature Reviews Neuroscience* 16, 7 (June 2015), 403–418. <https://doi.org/10.1038/nrn3976>
- [88] Stephen J Weber and Thomas D Cook. 1972. Subject effects in laboratory research: An examination of subject roles, demand characteristics, and valid inference. *Psychological Bulletin* 77, 4 (April 1972), 273–295. <https://doi.org/10.1037/h0032351>
- [89] John D. Wells, Damon E. Campbell, Joseph S. Valacich, and Mauricio Featherman. 2010. The Effect of Perceived Novelty on the Adoption of Information Technology Innovations: A Risk/Reward Perspective. *Decision Sciences* 41, 4 (Nov. 2010), 813–843. <https://doi.org/10.1111/j.1540-5915.2010.00292.x>
- [90] John R Wilson and Andrew Rutherford. 1989. Mental models: Theory and application in human factors. *Human Factors* 31, 6 (Dec. 1989), 617–634. <https://doi.org/10.1177/001872088903100601>
- [91] Arkady Zgonnikov, David Abbink, and Gustav Markkula. 2022. Should I Stay or Should I Go? Cognitive Modeling of Left-Turn Gap Acceptance Decisions in Human Drivers. *Human Factors* (Dec. 2022), 15 pages. <https://doi.org/10.1177/00187208221144561>
- [92] Rui Zhang, Nathan J. McNeese, Guo Freeman, and Geoff Musick. 2021. "An Ideal Human": Expectations of AI Teammates in Human-AI Teaming. In *Proceedings of the ACM on Human-Computer Interaction (CSCW3, Vol. 4)*. Article 246, 25 pages. <https://doi.org/10.1145/3432945>

A VERBAL DESCRIPTIONS OF THE SHAM-AI SYSTEM

The participants were presented with the following text as an introduction to the study. Depending on their group assignment (either DESCRIPTION), participants initially read the provided text. Followed by either a paragraph with a POSITIVE VERBAL DESCRIPTION or a NEGATIVE VERBAL DESCRIPTION, both concluding with the same paragraph.

The common paragraph was:

People perform more efficiently when the task difficulty level fits their stress level. Therefore, our team has developed ADAPTIMIND™, an AI system that adjusts task difficulty in reaction-critical contexts by analyzing the user's behavior and physiological signals, specifically the electrodermal activity (EDA) measured by medical-grade electrodes using two fingers of your hand.

Our AI system dynamically adjusts the task's difficulty by altering the task's pace according to your measured stress level. The algorithm is constantly learning from and adapting to the physiological indicators and your performance during the task. It may take some time to notice the changes in pace.

In the NEGATIVE DESCRIPTION condition, the following paragraph was then shown to participants:

The first users of ADAPTIMIND™ reported that when using the system, it decreased their task performance and increased stress making the task more difficult. As it is a new and untried AI system, it is very unreliable and risky to implement in real-world applications. In this study, we want to test these preliminary findings in a controlled setting.

For the POSITIVE DESCRIPTION condition the following paragraph was shown to the participants:

The first users of ADAPTIMIND™ reported that when using the system, it increased their task performance and decreased stress, making the task easier. As it is a cutting-edge AI system, it is very reliable and safe to implement in real-world applications. In this study, we want to test these preliminary findings in a controlled setting.

The text concluded in the same way for both groups:

We would like to evaluate your performance using AI and compare it to a condition where the AI is inactive (control condition). We will remind you in which of the two conditions you are in before starting the tasks.

B INFORMATION ON THE SHAM-AI SYSTEM STATUS - ACTIVE

Before the two blocks where participants performed the letter discrimination task with the sAI system active, the following text was displayed:

AI is now ACTIVE

The artificial intelligence system will now monitor your behavior and your physiological signals with the electrodes we have placed on your hand.

By monitoring your stress levels, the AI system will adjust the task pace. We will be assessing your performance based on reaction speed and accuracy.

The next paragraph differed based on the group allocation to positive/negative VERBAL DESCRIPTION:

POSITIVE VERBAL DESCRIPTION:

The system is expected to increase your task performance and decrease stress, making the task easier.

NEGATIVE VERBAL DESCRIPTION:

The system is expected to decrease your task performance and increase stress, making the task more difficult.

The text was concluded with the following instruction for both groups:

Please keep your hand with the electrodes on the table with your palm pointed upwards.

C INFORMATION ON THE SHAM-AI SYSTEM STATUS - INACTIVE

Before the two blocks where participants performed the letter discrimination task with the sAI system inactive, the following text was displayed:

AI is now INACTIVE

In this part of the study, we want to measure your performance without the AI system. The pace of the task will be random. We will be assessing your performance based on reaction speed and accuracy.

Please hold your hand on the table with your palm pointed upwards.

D MAILS, TIA AND SIPAS

Table 5: Mean scores and standard deviation as a function of DESCRIPTION for the questionnaires Meta AI Literacy Scale (MAILS), Checklist for Trust between People and Automation (TiA) and Subjective Information Processing Awareness Scale (SIPAS)

DESCRIPTION	MAILS		TiA		SIPAS	
	<i>M</i>	<i>SD</i>	<i>M</i>	<i>SD</i>	<i>M</i>	<i>SD</i>
Negative	108.61	28.28	47.97	9.15	3.47	1.05
Positive	118.71	27.89	47.39	9.93	3.58	1.03

E HIERARCHICAL DRIFT DIFFUSION MODEL WITH SYSTEM STATUS AND DESCRIPTION IN BRMS

All parameters are modeled on the log scale using the Wiener distribution.

Drift rate (v):

$$\begin{aligned} \log(v_{ijkl}) = & \beta_{0v} + \beta_{1v} \cdot \text{System Status}_j + \beta_{2v} \cdot \text{Description}_k \\ & + \beta_{3v} \cdot \text{System Status}_j \times \text{Description}_k \\ & + \beta_{4v} \cdot \text{Order}_l + b_{iv} \end{aligned}$$

Boundary separation (α):

$$\begin{aligned} \log(\alpha_{ijkl}) = & \beta_{0\alpha} + \beta_{1\alpha} \cdot \text{System Status}_j + \beta_{2\alpha} \cdot \text{Description}_k \\ & + \beta_{3\alpha} \cdot \text{System Status}_j \times \text{Description}_k \\ & + \beta_{4\alpha} \cdot \text{Order}_l + b_{i\alpha} \end{aligned}$$

Non-decision time (τ):

$$\begin{aligned} \log(\tau_{ijkl}) = & \beta_{0\tau} + \beta_{1\tau} \cdot \text{System Status}_j + \beta_{2\tau} \cdot \text{Description}_k \\ & + \beta_{3\tau} \cdot \text{System Status}_j \times \text{Description}_k \\ & + \beta_{4\tau} \cdot \text{Order}_l + b_{i\tau} \end{aligned}$$

Parameters and Priors:

Intercept priors:

$$\begin{aligned} \beta_{0v} & \sim \text{Normal}(0.74, 0.5) \\ \beta_{0\alpha} & \sim \text{Normal}(0.40, 1), \text{ lb} = 0.1 \\ \beta_{0\tau} & \sim \text{Normal}(-15, 1), \text{ lb} = -25, \text{ ub} = 3 \end{aligned}$$

Slope priors:

$$\begin{aligned} \beta_{1v}, \beta_{2v}, \beta_{3v}, \beta_{4v} & \sim \text{Normal}(0, 0.5) \\ \beta_{1\alpha}, \beta_{2\alpha}, \beta_{3\alpha}, \beta_{4\alpha} & \sim \text{Normal}(0, 0.1) \\ \beta_{1\tau}, \beta_{2\tau}, \beta_{3\tau}, \beta_{4\tau} & \sim \text{Normal}(0, 0.01) \end{aligned}$$

Random effects:

$$b_{iv}, b_{i\alpha}, b_{i\tau} \sim \text{Normal}(0, \sigma)$$

Reaction Time Modeling:

$$f(RT | \log(v_{ijkl}), \log(\alpha_{ijkl}), \log(\tau_{ijkl}), \text{bias} = 0.5)$$

Where RT is the observed reaction time.

F EMPIRICAL AND PREDICTED INDIVIDUAL REACTION TIME DIFFERENCE FOR SYSTEM STATUS

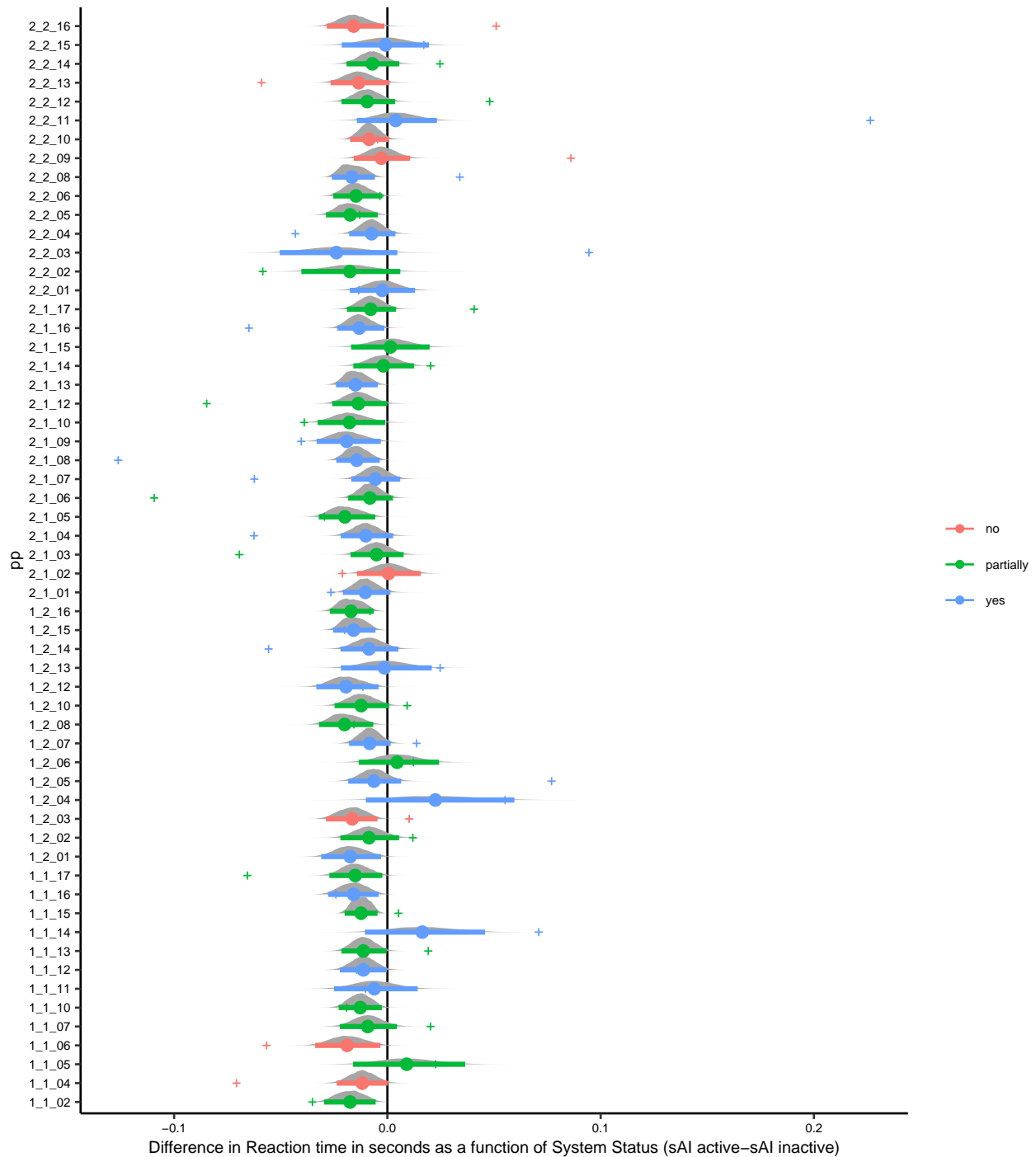


Figure 8: Individual difference (sAI active - sAI inactive) in reaction time predicted by the Drift-Diffusion model with 95% High-density intervals and the median estimate of the posterior distribution as a function of Manipulation Check (self-reported belief after the debriefing). + indicates the empirical mean difference in reaction time. Distance of empirical RT difference and predicted RT difference shows partial pooling as well as accounting for speed-accuracy trade-offs.

G DEVIATIONS FROM THE PRE-REGISTRATION

Table 6: Rationale for Deviations from Pre-Registration, for the pre-registration see <https://aspredicted.org/gm4n7.pdf>

Section	Deviation
Labels	We exchanged the term "nocebo" for the research questions and hypotheses with "negative verbal description" and "placebo" with "positive verbal descriptions." Additionally, the conditions were specified with "sham-AI" (sAI) in an active or inactive (control condition) state.
Participants: Recruiting and testing	We deviated from first testing 46 participants for nocebo (negative description), followed by testing 46 for placebo (positive description) due to time constraints. We stopped testing the negative description group after 30 participants were reached and then proceeded with testing the positive description group until we reached 60 participants. After this, we alternated the allocation of the last 6 participants to each group. The last day of testing remained the 18th of August 2023.
Performance data: Excluding trials	We excluded trials with too short responses by filtering RT under 150 ms instead of under 300 ms. This was a necessary deviation as participants were faster in their reactions than anticipated.
Performance data: Excluding participants	We excluded participants achieving less than 60% accuracy in any condition or having a miss rate exceeding 35%.
Performance data: Group Analyses	Given the AI performance bias, we modeled the data of both groups together instead of separately.

H MODEL PARAMETERS AND DIAGNOSTICS

H.1 Model found in Section 6.2.1

Table 7: Model Formula in Wilkinson notation: *Subjective overall performance rating* - 4 ~ 1 + *Description* × *System Status* + (1|*participant*)

Parameter	Median	p_b	\hat{R}	ESS	95% HDI	Prior
Intercept	0.51	0.00%	1.00	81228.44	[0.25, 0.77]	student (3, 0.50, 2.50)
Description	-0.19	6.70%	1.00	80719.48	[-0.45, 0.06]	normal (0, 1.39)
Time	0.19	4.66%	1.00	130714.43	[-0.03, 0.42]	normal (0, 1.39)
Description × Time	0.16	7.99%	1.00	136797.53	[-0.06, 0.38]	normal (0, 1.39)
$SD_{\text{participant}}$	0.50	0.00%	1.00	12043.44	[0.00, 0.85]	student (3, 0, 2.50)

Table 8: Model Formula in Wilkinson notation: *Subjective estimated task speed rating* - 50 ~ 1 + *Description* × *System Status* + (1|*participant*)

Parameter	Median	p_b	\hat{R}	ESS	95% HDI	Prior
Intercept	8.54	0.00	1.00	100738.47	[5.41, 11.78]	student (3, 7.50, 14.10)
Description	-1.31	21.09	1.00	103229.07	[-4.52, 1.88]	normal (0, 18.04)
Time	3.96	0.58	1.00	113865.65	[0.92, 7.03]	normal (0, 18.04)
Description×Time	-1.16	22.46	1.00	112055.75	[-4.17, 1.93]	normal (0, 18.04)
$SD_{\text{participant}}$	2.63	0.00	1.00	25022.39	[0, 6.90]	student (3, 0, 14.10)

Table 9: Model Formula in Wilkinson notation: *Subjective estimated number of correct responses* $\sim 1 + \text{Description} \times \text{System Status} \times \text{Time} + (1|\text{participant})$

Parameter	Median	p_b	\hat{R}	ESS	95% HDI	Prior
Intercept	136.30	0.00	1.00	21231.73	[129.11, 143.47]	student (3, 150, 44.50)
Description	-1.38	35.22	1.00	20697.38	[-8.68, 5.75]	normal (0, 36.71)
System Status	6.32	0.00	1.00	125110.60	[3.29, 9.30]	normal (0, 36.71)
Time	6.31	0.00	1.00	130118.28	[3.36, 9.28]	normal (0, 36.71)
Description×System Status	0.05	48.66	1.00	125738.08	[-2.97, 3.07]	normal (0, 36.71)
Description×Time	-2.75	3.74	1.00	131880.02	[-5.76, 0.29]	normal (0, 36.71)
System Status×Time	2.63	4.46	1.00	131160.58	[-0.39, 5.62]	normal (0, 36.71)
Description×System Status×Time	-0.07	48.21	1.00	131327.49	[-3.18, 2.90]	normal (0, 36.71)
$SD_{\text{participant}}$	26.56	0.00	1.00	20598.96	[21.06, 32.61]	student (3, 0, 44.50)

Table 10: Model Formula in Wilkinson notation: *Reaction time (s) * 1000* $\sim \text{System Status} + \text{Description} + \text{Order} + (1|\text{participant})$

Parameter	Median	p_b	\hat{R}	ESS	95% HDI	Prior
Intercept	606.97	0.00	1.00	2927.02	[585.57, 628.2]	student (3, 567.10, 123.10)
System Status	-4.17	0.00	1.00	28804.59	[-6.18, -2.19]	normal (0, 171.61)
Description	-5.00	32.72	1.00	2501.51	[-25.79, 16.94]	normal (0, 171.61)
Order	11.06	0.00	1.00	30137.34	[9.03, 13.01]	normal (0, 171.61)
$SD_{\text{participant}}$	80.75	0.00	1.00	3624.85	[66.9, 97.36]	student (3, 0, 123.10)

Table 11: Model Formula in Wilkinson notation: *Correctness of responses* $\sim \text{System Status} + \text{Description} + \text{Order} + (1|\text{participant})$

Parameter	Median	p_b	\hat{R}	ESS	95% HDI	Prior
Intercept	2.41	0.00	1.00	6588.72	[2.19, 2.64]	student (3, 0, 2.50)
System Status	0.05	2.05	1.00	61408.52	[0, 0.09]	student (3, 0, 10)
Description	0.10	19.00	1.00	6137.11	[-0.12, 0.33]	student (3, 0, 10)
Order	-0.05	1.37	1.00	58438.11	[-0.09, 0]	student (3, 0, 10)
$SD_{\text{participant}}$	0.83	0.00	1.00	8030.09	[0.68, 1.01]	student (3, 0, 2.50)

Table 12: Model Formula in Wilkinson notation: $RTsec|dec(lu) \sim \text{System Status} \times \text{Description} + \text{Order} + (1|p|participant)$, Model outputs for the parameters on the log scale. Medians are provided for each parameter, along with their 95% HDI and p_b . Parameters distinguishable from zero are marked with *. We ran the model with two chains and 4000 iterations.

Parameter	Median	p_b	\hat{R}	ESS	95% HDI	Prior
ν Intercept	0.68	0.00	1.00	1097.76	[0.55, 0.81]	normal (0.74, 0.50)
α Intercept	0.37	0.00	1.00	1700.25	[0.32, 0.43]	normal (0.40, 1)
τ Intercept	-1.21	0.00	1.00	1157.78	[-1.29, -1.13]	normal (-15, 1)
ν System Status	0.03	0.00	1.00	10716.08	[0.02, 0.05]	normal (0, 0.50)
ν Description	0.06	18.11	1.00	1308.32	[-0.06, 0.19]	normal (0, 0.50)
ν Order	-0.02	0.30	1.00	10441.23	[-0.03, 0]	normal (0, 0.50)
ν System Status×Description	0.01	7.68	1.00	11569.49	[0, 0.02]	normal (0, 0.50)
α System Status	0.01	0.25	1.00	10301.47	[0, 0.02]	normal (0, 0.10)
α Description	0.02	13.90	1.00	1539.17	[-0.02, 0.07]	normal (0, 0.10)
α Order	0.01	14.23	1.00	9481.10	[0, 0.01]	normal (0, 0.10)
α System Status×Description	0.01	9.55	1.00	10416.62	[0, 0.02]	normal (0, 0.10)
τ System Status	-0.02	0.00	1.00	13135.76	[-0.02, -0.02]	normal (0, 0.01)
τ Description	0.00	49.66	1.00	7734.16	[-0.02, 0.02]	normal (0, 0.01)
τ Order	0.01	0.00	1.00	13118.79	[0.01, 0.02]	normal (0, 0.01)
τ System Status×Description	-0.01	0.39	1.00	12630.83	[-0.01, 0]	normal (0, 0.01)
$SD_{participant}$	0.49	0.00	1.00	2011.08	[0.40, 0.60]	student (3, 0, 2.50)
$SD_{participant} \times \alpha$ Intercept	0.20	0.00	1.00	2108.20	[0.16, 0.24]	student (3, 0, 2.50)
$SD_{participant} \times \tau$ Intercept	0.29	0.00	1.00	1954.95	[0.23, 0.34]	student (3, 0, 2.50)
$cor_{participant} \times \nu$ Intercept× α Intercept	0.11	20.39	1.00	2067.49	[-0.15, 0.36]	lkj(2)
$cor_{participant} \times \nu$ Intercept× τ Intercept	0.31	0.98	1.00	1813.20	[0.06, 0.54]	lkj(2)
$cor_{participant} \times \alpha$ Intercept× τ Intercept	-0.56	0.00	1.00	2043.48	[-0.73, -0.36]	lkj(2)

Table 13: Model Formula in Wilkinson notation: Cognitive Workload (TLX_{sum}) $\sim \text{Description} \times \text{System Status} + \text{Order} + (1|participant)$

Parameter	Median	p_b	\hat{R}	ESS	95% HDI	Prior
Intercept	64.75	0.00	1.00	23700.70	[60.28, 69.33]	student (3,64, 19.30)
Description	0.83	35.91	1.00	22840.38	[-3.65, 5.50]	normal (0, 20.34)
System Status	-0.08	46.82	1.00	136110.80	[-2.03, 1.91]	normal (0, 20.34)
Order	1.45	7.28	1.00	137885.40	[-0.57, 3.37]	normal (0, 20.34)
Description×System Status	0.43	33.28	1.00	140187.50	[-1.55, 2.38]	normal (0, 20.34)
$SD_{participant}$	16.69	0.00	1.00	19669.48	[13.16, 20.54]	student (3, 0, 1)

Table 14: Model Formula in Wilkinson notation: Physiological Arousal (SCL_{mean}) $\sim \text{System Status} \times \text{Description} + \text{Order} + (1|participant)$

Parameter	Median	p_b	\hat{R}	ESS	95% HDI	Prior
Intercept	-0.96	0.00	1.00	72495.54	[-1.32, -0.59]	student (3, -0.80, 2.50)
System Status	-0.20	11.88	1.00	48124.82	[-0.53, 0.13]	normal (0, 0.90)
Description	0.25	7.37	1.00	71529.11	[-0.10, 0.59]	normal (0, 0.90)
Order	0.20	4.20	1.00	75567.21	[-0.03, 0.42]	normal (0, 0.90)
System Status×Description	0.07	32.70	1.00	52413.58	[-0.24, 0.38]	normal (0, 0.90)
$SD_{participant}$	0.15	0.00	1.00	21592.19	[0, 0.36]	student (3, 0, 2.50)

Table 15: Model Formula in Wilkinson notation: System evaluation item 1 – 4 ~ Description

Parameter	Median	p_b	\hat{R}	ESS	95% HDI	Prior
Intercept	-0.60	0.45	1.00	64497.24	[-1.03, -0.16]	student (3, -1, 2.50)
Description	-0.17	21.58	1.00	65782.05	[-0.63, 0.25]	normal (0, 1.72)

Table 16: Model Formula in Wilkinson notation: System evaluation item 2 – 4 ~ Description

Parameter	Median	p_b	\hat{R}	ESS	95% HDI	Prior
Intercept	-0.24	12.54	1.00	64462.65	[-0.66, 0.16]	student (3, 0, 2.50)
Description	-0.41	2.48	1.00	65215.93	[-0.82, 0]	normal (0, 1.66)

Table 17: Model Formula in Wilkinson notation: System evaluation item 3 – 4 ~ Description

Parameter	Median	p_b	\hat{R}	ESS	95% HDI	Prior
Intercept	-0.32	6.39	1.00	68372.93	[-0.74, 0.10]	student (3, 0, 2.50)
Description	-0.35	5.15	1.00	65799.87	[-0.76, 0.07]	normal (0, 1.67)

Table 18: Model Formula in Wilkinson notation: System evaluation item 4 – 4 ~ Description

Parameter	Median	p_b	\hat{R}	ESS	95% HDI	Prior
Intercept	-0.33	5.87	1.00	66744.43	[-0.74, 0.08]	student (3, 0, 2.50)
Description	-0.27	9.35	1.00	67963.57	[-0.68, 0.14]	normal (0, 1.62)

Table 19: Model Formula in Wilkinson notation: System evaluation item 5 – 4 ~ Description

Parameter	Median	p_b	\hat{R}	ESS	95% HDI	Prior
Intercept	0.02	44.47	1.00	59849.21	[-0.33, 0.39]	student (3, 0, 2.50)
Description	-0.12	24.80	1.00	61742.01	[-0.48, 0.23]	normal (0, 1.42)

Table 20: Model Formula in Wilkinson notation: System evaluation item 6 – 4 ~ Description

Parameter	Median	p_b	\hat{R}	ESS	95% HDI	Prior
Intercept	0.09	31.37	1.00	64222.85	[-0.27, 0.45]	student (3, 0, 2.50)
Description	-0.11	28.33	1.00	63489.53	[-0.46, 0.27]	normal (0, 1.44)

Table 21: Model Formula in Wilkinson notation: System evaluation item 7 – 4 ~ Description

Parameter	Median	p_b	\hat{R}	ESS	95% HDI	Prior
Intercept	-0.35	2.06	1.00	61105.53	[-0.68, -0.01]	student (3, 0, 2.50)
Description	-0.09	29.47	1.00	59780.62	[-0.43, 0.25]	normal (0, 1.36)

Table 22: Model Formula in Wilkinson notation: System evaluation item 8 – 4 ~ Description

Parameter	Median	p_b	\hat{R}	ESS	95% HDI	Prior
Intercept	1.23	0.00	1.00	66607.93	[0.93, 1.51]	student (3, 1, 2.50)
Description	-0.20	8.98	1.00	67891.85	[-0.49, 0.09]	normal (0, 1.15)

Table 23: Model Formula in Wilkinson notation: UEQ-S-pragmatic – 0 ~ Description

Parameter	Median	p_b	\hat{R}	ESS	95% HDI	Prior
Intercept	0.73	0.00	1.00	57789.76	[0.46, 0.99]	student (3, 0.80, 2.50)
Description	-0.17	10.79	1.00	59970.21	[-0.43, 0.10]	normal (0, 1.06)

Table 24: Model Formula in Wilkinson notation: UEQ-S-hedonic – 0 ~ Description

Parameter	Median	p_b	\hat{R}	ESS	95% HDI	Prior
Intercept	0.78	0.00	1.00	65655.38	[0.49, 1.08]	student (3, 1, 2.50)
Description	-0.04	40.80	1.00	68065.81	[-0.33, 0.26]	normal (0, 1.17)

Table 25: Model Formula in Wilkinson notation: SUS-Adapted Score – 68 ~ Description

Parameter	Median	p_b	\hat{R}	ESS	95% HDI	Prior
Intercept	-1.41	23.55	1.00	65479.48	[-5.39, 2.47]	student (3, -0.50, 14.80)
Description	-1.71	19.28	1.00	65245.98	[-5.69, 2.18]	normal (0, 15.65)

Table 26: Replication Study - Model Formula in Wilkinson notation: Expected overall performance – 4 ~ 1 + Description

Parameter	Median	p_b	\hat{R}	ESS	95% HDI	Prior
Intercept	0.85	0.00	1.00	125379.81	[0.63, 1.08]	student (3, 1, 2.50)
Comprehension	-0.26	1.05	1.00	122476.12	[-0.48, -0.04]	normal (0, 1.08)

Table 27: Replication Study - Model Formula in Wilkinson notation: Expected task speed – 50 ~ 1 + Description

Parameter	Median	p_b	\hat{R}	ESS	95% HDI	Prior
Intercept	15.37	0.00	1.00	65594.30	[11.73, 19.11]	student (3, 15, 19.30)
Comprehension	-4.60	0.82	1.00	61874.58	[-8.34, -0.86]	normal (0, 18.34)

Table 28: Replication Study - Model Formula in Wilkinson notation: Estimated correct ~ 1 + Comprehension × System Status + (1|participant)

Parameter	Median	p_b	\hat{R}	ESS	HDI	Prior
Intercept	147.23	0.00	1.00	35543.08	[139.7, 154.31]	student (3, 150, 44.50)
Comprehension	-0.51	44.50	1.00	38871.60	[-7.80, 6.74]	normal (0, 39.43)
System Status	5.71	0.15	1.00	151821.25	[2, 9.39]	normal (0, 39.43)
Comprehension×System Status	-4.36	1.03	1.00	146724.78	[-8.08, -0.69]	normal (0, 39.43)
$SD_{\text{participant}}$	30.22	0.00	1.00	20042.75	[24.22, 37]	student (3, 0, 44.5)