# Diverse Perspectives on AI: Examining People's Acceptability and Reasoning of Possible AI Use Cases

JIMIN MUN, Carnegie Mellon University, USA

WEI BIN AU YEONG, Carnegie Mellon University, USA

WESLEY HANWEN DENG, Carnegie Mellon University, USA

JANA SCHAICH BORG, Duke University, USA

MAARTEN SAP, Carnegie Mellon University, USA

In recent years, there has been a growing recognition of the need to incorporate lay-people's input into the governance and acceptability assessment of AI usage. However, how and why people judge different AI use cases to be acceptable or unacceptable remains under-explored. In this work, we investigate the attitudes and reasons that influence people's judgments about AI's development via a survey administered to demographically diverse participants (N=197). We focus on ten distinct professional (e.g., Lawyer AI) and personal (e.g., Digital Medical Advice AI) AI use cases to understand how characteristics of the use cases and the participants' demographics affect acceptability. We explore the relationships between participants' judgments and their rationales such as reasoning approaches (cost-benefit reasoning vs. rule-based). Our empirical findings reveal number of factors that influence acceptance such as general negative acceptance and higher disagreement of professional usage over personal, significant influence of demographics factors such as gender, employment, and education as well as AI literacy level, and reasoning patterns such as rule-based reasoning being used more when use case is unacceptable. Based on these findings, we discuss the key implications for soliciting acceptability and reasoning of AI use cases to collaboratively build consensus. Finally, we shed light on how future FAccT researchers and practitioners can better incorporate diverse perspectives from lay people to better develop AI that aligns with public expectations and needs.

## 1 Introduction

There are growing calls to regulate AI's development and integration into society [60]. These efforts, as reflected in the EU AI Act [59], NIST AI Risk Management framework [55], and recent U.S. Executive Order [1], have resulted in discussions about whether certain AI use cases should be pursued at all. Despite much progress in this area, it is still not clear how to determine which use cases should be pursued or more heavily regulated. Further, little is known about how lay community members, especially those from marginalized groups, feel about the development of specific AI use cases [2, 71].

One significant challenge when evaluating the acceptability and impact of specific AI use cases[1] is that there can be both positive and negative effects, depending on the context [53]. For instance, while educational AI can provide affordable and accessible personal tutor, it can also lead to over-reliance of students and diminish the goal of education [75, 83]. To develop a generalizable approach to making decisions about AI use cases, ideally we would understand how people resolve these conflicts. More specifically, first, it is essential to understand differences in judgments about *acceptability and likely usage* across use cases, and how such judgments relate to scenario characteristics. Second, we need better understanding of the *personal factors influencing these judgments*, especially as they relate to demographic differences [42]. Third, we need to better understand the reasoning strategies participants use when making judgments about AI use cases, and how those strategies do or do not relate to the judgments that are ultimately made. To form

---

[1]By use cases, we mean specific scenarios, applications, or problems that an AI system is designed to solve or assist within real-world contexts.

Authors' Contact Information: Jimin Mun, jmun@andrew.cmu.edu, Carnegie Mellon University, Pittsburgh, Pennsylvania, USA; Wei Bin Au Yeong, wauyeong@andrew.cmu.edu, Carnegie Mellon University, Pittsburgh, Pennsylvania, USA; Wesley Hanwen Deng, hanwend@andrew.cmu.edu, Carnegie Mellon University, Pittsburgh, Pennsylvania, USA; Jana Schaich Borg, js524@duke.edu, Duke University, Durham, North Carolina, USA; Maarten Sap, msap2@andrew.cmu.edu, Carnegie Mellon University, Pittsburgh, Pennsylvania, USA.

governance and policy decisions that anticipate and address disagreements about the development or regulation about specific AI use cases, these understandings are crucial, especially across groups of people with diverse backgrounds, experiences, and familiarity with AI.

To address these needs, in this work we examine how and why lay people judge various AI use cases as acceptable or unacceptable, asking the following research questions:

**RQ1** How do judgments of acceptability vary across a set of distinct AI use cases and their characteristics?

**RQ2** What attributes or characteristics of people explain the variation in acceptability judgments?

**RQ3** How do people reason through acceptability judgments of AI use cases?

To answer these questions, we develop a survey to collect judgments and reasoning processes of 197 demographically diverse participants with varying levels of experience with AI. We ask participants to report whether a certain AI use case should be developed or not, whether they would use such a system, and ask them to provide rationales for their judgment and conditions that would cause them to change their judgments (Figure 1). We examine ten different AI use cases[2]. To account for differences between sectors or domains, we select use cases in two categories, professional and personal use, and vary them by required entry-level education and EU AI risk level (Table 1).

To meaningfully differentiate and analyze participants' reasons and reasoning strategies, we borrow concepts from moral psychology and philosophy. We investigate participants' rationales through two distinct but sometimes overlapping reasoning patterns: cost-benefit reasoning, which assesses expected outcomes (e.g., "using AI for this task would save time"), and rule-based reasoning, which evaluates the intrinsic values of the action itself (e.g., "having humans perform this task would be inherently wrong") [21, 22]. We further explore the moral foundations reflected in participants' reasoning, with moral foundations theory[3] [33, 34]. Additionally, to understand aspects of AI that raise concerns, we employ three dimensions based on prior studies [53, 68]: functionality (system capabilities like performance, bias, and privacy), usage (context of system integration, such as supervision, misuse, or unintended use), and societal impact (effects on individuals, communities, and society, such as job loss and over-reliance).

Our empirical results show general higher acceptance of personal use cases over professional. While both categories of use cases show decreased acceptance with increased entry level education and risk, professional use cases display more variability and disagreements across judgments (*RQ1*). Acceptability significantly varied among demographic groups and levels of AI literacy, with lower acceptability observed particularly among non-male participants and those familiar with AI ethics (RQ2). Finally, our results show varying distribution of reasoning types across acceptability decisions with rule-based reasoning being associated with negative acceptance as well as concern for societal impact. Further qualitative analysis reveal rules such as the need for humanness in certain use cases whether it be for empathy or interaction (*RQ3*).

Our findings shed novel light onto the diversity of people's acceptability and reasoning of AI uses in distinct domains and risk levels. We conclude with a discussion highlighting three key implications: first, diverse methodologies are needed to effectively analyze use cases and their characteristics; second, involving diverse stakeholders is crucial for assessing the acceptability of AI applications, particularly in workplaces; and third, further investigation into human reasoning processes about AI, notably rule-based reasoning, is needed to inform consensus-building in policy making.

---

[2]We focus on text-based, non-embodied, digital systems, and while we do not specifically discuss the AI user and subject, in our use case description, we follow three of the five concepts used in EU AI Act to describe high risk use cases [31]: the domain, purpose, and capabilities.

[3]We used the five foundational dimensions: Care, Fairness, Loyalty, Authority, and Purity. Although these dimensions have been updated to encompass a broader range of values beyond WEIRD (White, Educated, Industrialized, Rich, and Democratic) populations [7], we selected this version for survey brevity.

## 2  Related Works

In this section, we briefly summarize the background and related work towards assessing acceptability and impact of AI use cases. In each subsection, we highlight how our work extends prior work.

### 2.1  Assessing Impact of AI

Recent years have witnessed increasing calls from academics [3, 6, 11, 17, 23, 35, 37, 40, 47, 56], government [54, 73, 74], civil society [4, 48, 57, 62, 66, 73], and industry [24, 35, 50, 51, 58, 76] to assess the impact of AI systems designed and developed by AI researchers and practitioners. This effort has particularly highlighted the need to understand the *positive* impact while grappling with the potential *negative* impact of integrating AI into certain products and services that affect people's daily lives [6, 35, 47]. In response, researchers in FAccT, HCI, and AI have developed tools and processes to support AI researchers and practitioners in anticipating the impact of AI systems they developed[16, 24, 40, 78, 79]. For example, many have developed AI impact taxonomies or checklists to help developers categorize AI impact[51, 79]. Wang et al. and Deng et al. developed tools and templates to support industry AI developers and researchers in assessing the potential negative societal impact of their work, such as job displacement or stereotyping social groups.

However, this prior work primarily focuses on supporting *AI experts* rather than *diverse lay people*'s impact assessments of potential AI use cases. Our work extends these prior efforts by understanding diverse (and sometimes conflicting) perspectives on both positive and negative impact of AI use cases from lay people, as a crucial step to complement the AI impact assessments conducted by AI researchers and practitioners.

### 2.2  Understanding People's Perceptions of AI Use Cases

Responding to the calls on meaningfully engaging lay people in assessing the impact of specific AI use cases, prior work have started to understand lay people's perceptions of AI use case [16, 40, 42, 53]. Among other findings, this prior work revealed a substantial amount of disagreement regarding the desired behavior of AI, primarily due to the subjectivity inherent in certain tasks (e.g., toxicity detection [14, 63], image captioning [85]) and ambiguous ethical implication of decisions made by AI for certain tasks (e.g., self driving cars [8], medical AI [18], predictive analysis [10]). Work done by Mun et al. highlighted that lay people can envision diverse set of harms specific to different AI use cases, complementary to those defined by experts. Another line of work also begins to examine how factors such as demographic backgrounds and previous exposure to discrimination can affect people's sensitivity towards potential AI harms [42].

Our work extends this prior work by examined the **detailed reasoning processes** of lay people regarding the acceptability and the **trade-offs between positive and negative impacts** of AI use cases. In particular, we draw on model decision theory framework, such as the moral foundations developed by Graham et al., to design a survey flow (See Figure 1 to solicit lay people's decision-making processes (and potential moral conflicts) when assessing both the benefits and harms of concrete AI use cases.

### 2.3  Background: Moral Decision Making

Morality, characterized by diverse values across cultures and social groups, aims to suppress selfishness to facilitate social life [39]. To understand decision-making in AI use cases, we draw on moral psychology and dual system theory. We examine two decision-making systems: cost-benefit reasoning, which assesses outcomes and consequences, and rule-based reasoning, focusing on norms, rules, and virtues [21, 22]. These correspond to utilitarian reasoning (maximizing

good) and deontological reasoning (duties and rights), respectively. Additionally, we apply moral foundations theory [34] to identify values and potential moral conflicts in AI development.
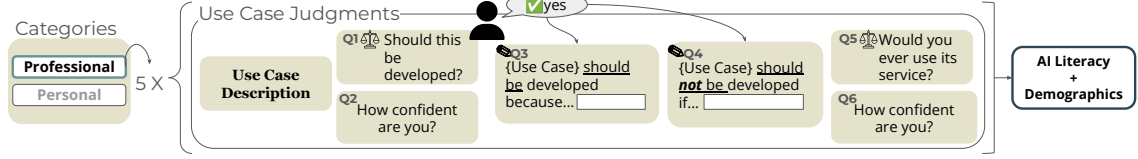


Fig. 1. Five professional or personal use cases are presented in a random order. For each use case, we ask multiple-choice questions about its development and confidence levels (Q1, Q2), free-text questions on rationale and decision-switching conditions (Q3, Q4), and multiple-choice questions on usage and confidence (Q5, Q6). These are followed by questions on AI literacy and demographics.

## 3 Study Design and Data Collection

To understand the decision making factors and reasoning patterns of a diverse population regarding AI use cases, we conducted a survey-based study with demographically diverse participants. In this section, we discuss the selection of use cases (§ 3.1), survey design (§ 3.2), and data collection details and participant demographics (§ 3.3).

### 3.1 Use Cases

To understand how different characteristics of AI use cases can impact judgments and decision making processes, we crafted ten different AI use cases. We first chose two broad application categories frequently mentioned by the public in previous works [41, 53]: **AI in personal**, everyday usage where participants could uniformly consider themselves as AI users and **AI in professional usage** where AI takes on a role thus far done by a human as a profession, e.g., Lawyer AI. We then developed five use cases in the AI in professional usage category that varied according to the level of education required for entry and five use cases in AI in personal category that varied according to the risk level assigned by the EU AI act. All the personal use cases were in the health domain, reflecting a key domain of interest conveyed by lay users in previous works [41, 53].

| Use Case | Factor | Description |
|---|---|---|
| **Professional Use Cases** | | |
| Lawyer | Doctoral/Professional Degree | Advises clients on digital legal proceedings/transactions. |
| Elementary School Teacher | Bachelor's degree | Teaches academic skills at the elementary school level. |
| IT Support Specialist | Some college, no degree | Maintains computer networks and provides technical help. |
| Government Eligibility Interviewer | High school diploma | Determine eligibility for government programs/resources. |
| Telemarketer | No formal education | Solicits donations or orders over the telephone. |
| **Personal Use Cases** | | |
| Digital Medical Advice | High Risk | Provide medical assessments prior to medical consultations. |
| Customized Lifestyle Coach | High / Limited Risk | Personalized advice for healthy living and wellness. |
| Personal Health Research | Limited Risk | Summarizes research related to personal health issues. |
| Nutrition Optimizer | Limited / Low Risk | Personalize meals and optimize nutritional intake. |
| Flavorful Swaps | Low Risk | Suggest delicious and healthy alternatives food options. |

Table 1. Use cases selected for our study by categories. Use case descriptions were shortened for brevity.

*Professional Use Case Scenarios.* For the first area of focus, AI in labor replacement, we collected jobs listed in the U.S. census bureau[4] and sorted them according to entry level education required as stated in the census. We chose education level as it has been tightly linked to socioeconomic and occupational status [29, 72]. We selected jobs that have a large portion of digital or intellectual components with minimal requirement for embodiment resulting in following five professional roles: Lawyer, Elementary school teacher, IT support specialist, Government support eligibility interviewer, and Telemarketer. See Table 1 for further details.

*Personal Use Case Scenarios.* To understand the acceptability of different health applications in personal and private life, we adapted the descriptions of personal use cases written by participants from prior works [41, 53]. The research team iteratively refined the descriptions to reflect risk levels according to EU AI Act, and confirmed agreement with categories assigned by GPT-4, following Herdel et al.. See Table 1 for further details.

## 3.2 Survey Design

Our survey presents participants with five use case descriptions, all from either the professional or personal category (randomly assigned and presented in random order) (see § 3.1 for details). After each description, participants answer: "Do you think a technology like this should be developed?" (Q1) and then, "How confident are you in your above answer?" (Q2). They provide open-text rationales by completing the prompt, "[Use Case] should be developed because..." (Q3), tailored to their Q1 response ("should" for "Yes" and "should not" for "No"). Participants also specify a condition for switching their decision with, "[Use Case] should not be developed if..." (Q4), adjusted similarly to Q1. Subsequently, they answer, "If [Use Case] existed, would you use its service?" (Q5) and express confidence with, "How confident are you in your above answer?" (Q6). Refer to Table 6 in the Appendix for the exact wording of the questions.

*Collecting Participant Characteristics.* Following the main survey, we asked participants questions about their AI literacy and demographics to explore various factors affecting perception of AI acceptance. We adopted a shortened version of AI literacy questionnaires from previous works [53, 77] with four AI literacy aspects, AI awareness, usage, evaluation, and ethics, and two additional questions for generative AI, usage frequency and familiarity with limitations. We collected demographic information of the participants such as race, gender, age, sexual orientation, religion, employment status, income, and level of education. Additionally, we collected information about chronicity, i.e., prolonged experiences of everyday discrimination, of their discrimination experiences (if any) following Kingsley et al..

## 3.3 Data Collection and Participant Demographics

We used Prolific[5] to recruit participants. To represent diverse sample, we stratified our recruitment by the ethnicity categories (White, Mixed, Asian, Black and Other) and age (18-48, 49-100) provided by Prolific. We also added criteria for quality such as survey approval rating and number of previous surveys completed. Our study was approved by IRB at our institutions, and we paid 12 USD/hour. Our final sample consisted of 197 participants across two categories, with professional usage assigned to 100 participants and personal to 97. See Appendix A.2 for further details on participants.

## 4 Analysis Methods

Our surveys consisted of both multiple choice (numerical) and open-text questions designed to answer our research questions. In this section, we detail our process for numerical (§ 4.1) and open-text (§ 4.2) analysis.

---

[4]https://www.bls.gov/ooh/occupation-finder.htm
[5]https://www.prolific.com

### 4.1 Multiple Choice Analysis

We analyzed the judgment and confidence ratings by mapping judgment (Q1, Q5) to 1 ("Should be developed", "Would use") or -1 ("Should not be developed", "Would not use") and confidence (Q2, Q6) to a scale from 1 to 5. We used numerically converted judgment, confidence, and combined (judgment×confidence; -5 to 5) values as dependent variables in our analysis. We used repeated-measures ANOVAs to understand the differences in mean responses between conditions/groups and linear mixed effects regression models (lmer) to better understand the effects of specific factors. We included a subject-specific random effect when using ANOVA and regression models and added a use-case-specific random effect when applicable. We factorized demographic responses for analysis with the exception of discrimination chronicity, which we aggregated to a numerical value [42, 49]. We also converted responses to AI literacy questions to numerical values for analysis.

### 4.2 Open-response Analysis

To assess the reasoning methods used by the participants, we analyzed the open-text responses on elaborations to their decisions (Q3) and circumstances in which their decisions would switch (Q4) along the following three dimensions: reasoning types (cost-benefit, rule-based, both, unclear), reference to moral foundations (Care, Fairness, Purity, Authority, Loyalty), and switching conditions (Functionality, Usage, Societal Impact). By analyzing reasoning types and moral values reflected in the participants' justifications, we aim to characterize *how* participants made their decisions, and by analyzing various factors such as primary concerns and stakeholders, we aim to discover *what* aspects were salient for the participants in their decisions.

*Classification and Aggregation.* We classified participants' responses to Q3 (elaboration of judgment) and Q4 (conditions for switching decisions) using OpenAI's o1-mini[6]. To validate the model's classification performance, results were compared with a reference set of 100 samples annotated by three independent annotators, comprised of team members and a professional annotator. Initially, each annotator independently assessed the data, and then consensus was reached through discussion to establish a gold standard set. The inter-rater agreement between the gold standard and o1-mini's annotations was evaluated using Gwet's AC1 metric, chosen for its robustness with infrequent labels [81]. The agreement levels varied, ranging from moderate to almost perfect with no dimension below moderate. Annotations for cost-benefit reasoning, rule-based reasoning, and authority reached near-perfect agreement; purity and usage achieved substantial agreement; the rest showed moderate agreement. Due to a lack of sufficient test samples and minimal occurrences in the annotated data, the moral value dimension Loyalty was excluded from further analysis. The annotations were converted into a one-hot encoding format for statistical analysis. See Appendix B for further details.

## 5 Findings

Our work aims to uncover variations in acceptability of AI use cases and factors and reasoning processes that underlie these judgments. In this section, we discuss our findings about the judgments of the AI use cases (§5.1), personal factors that may influence the decision such as demographics and AI literacy (§5.2), and factors in rationales that could uncover reasoning processes that lead to judgments (§5.3).

### 5.1 RQ1. Perceptions and Disagreements of Use Cases
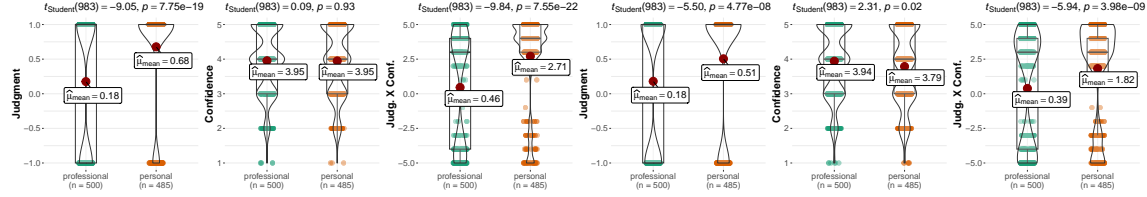
---

[6]`o1-mini-2024-09-12`

Fig. 2. Use case category means and distributions of numerically converted Judgment, Confidence, and Judgment×Confidence. Left three show results for development decisions (Q1, Q2) and right three show results for usage decisions (Q5, Q6). Significance was calculated using Student's t-test as indicated by labels above each plot.

*5.1.1 Use Case Factors.* In our analysis, we investigated the effects of use cases on participants' judgments of two categories—professional and personal—along with ten specific use cases on judgments and confidence about their development and usage. Figure 2 illustrates significant differences in acceptability depending on the categories. Notably, personal use cases had higher acceptability ($M_{DEV} = 0.68, SD_{DEV} = 0.74; M_{USAGE} = 0.51, SD_{USAGE} = 0.86$) than professional use cases ($M_{DEV} = 0.18, SD_{DEV} = 0.99; M_{USAGE} = 0.18, SD_{USAGE} = 0.98$). These differences are more pronounced when judgments are weighted by confidence. Although levels of confidence did not differ significantly between categories, they were fairly high for both categories, and personal category exhibited slightly lower confidence for usage.



Fig. 3. Professional use case means and distributions of numerically converted Judgment, Confidence, and Judgment×Confidence. First row shows results for development decisions (Q1, Q2) and second row shows results for usage decisions (Q5, Q6). ANOVA results for use cases are shown above each panel. Within subject test was performed using Student's t-test with Holm correction. * denotes following significant p-values: *** $p < 0.001$; ** $p < 0.01$; * $p < 0.05$. Use case names were shortened.

*Professional Use Cases.* Exploring specific use cases within the professional category, we observed that Elementary School Teacher AI ($M_{DEV} = -0.24, SD_{DEV} = 0.98; M_{USAGE} = -0.14, SD_{USAGE} = 1.00$) had the lowest acceptability for both types of judgments followed by Lawyer AI ($M_{DEV} = 0.04, SD_{DEV} = 1.00$) and Telemarketer AI ($M_{USAGE} = -0.08, SD_{USAGE} = 1.00$). Interestingly, IT Support Specialist AI had the highest acceptability ($M_{DEV} = 0.66, SD_{DEV} = 0.76; M_{USAGE} = 0.78, SD_{USAGE} = 0.63$). While all other use cases showed higher acceptability for usage over development, Telemarketer AI uniquely had higher acceptance for development over usage ($M_{DEV} = 0.26, SD_{DEV} = 0.97; M_{USAGE} = -0.08, SD_{USAGE} = $

1.00). Additionally, confidence on using Telemarketer AI ($M = 4.14, SD = 0.93$) was significantly ($p < .05$) higher than that of Government Eligibility Interviewer AI ($M = 3.79, SD = 0.99$), which had the lowest confidence in usage. Near 0 mean for development of Lawyer AI (0.04) and usage for Telemarketer AI ($-0.08$) suggest disagreement within judgments. Moreover, Elementary School Teacher AI uniquely unacceptable across both acceptability judgments ($-0.24$, $-0.14$) underscoring a unique characteristic perhaps related to care. See Figure 3 for further details
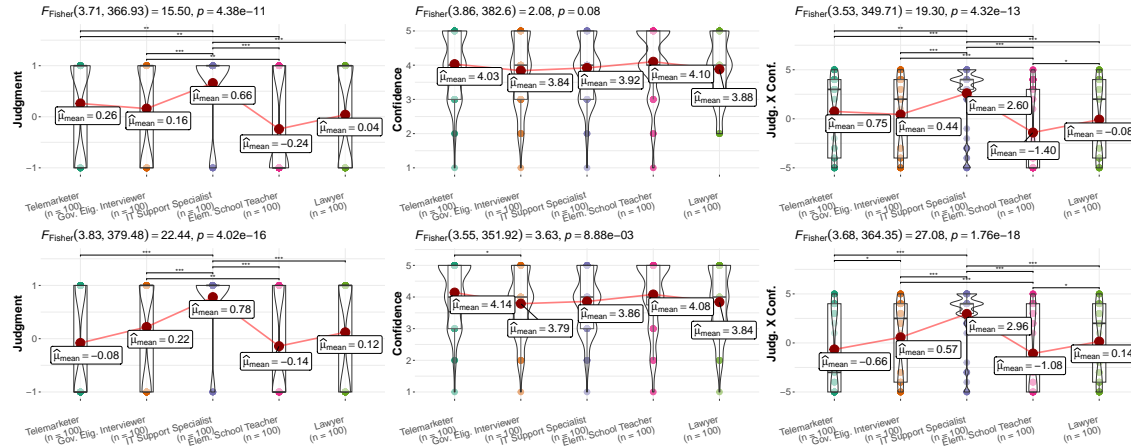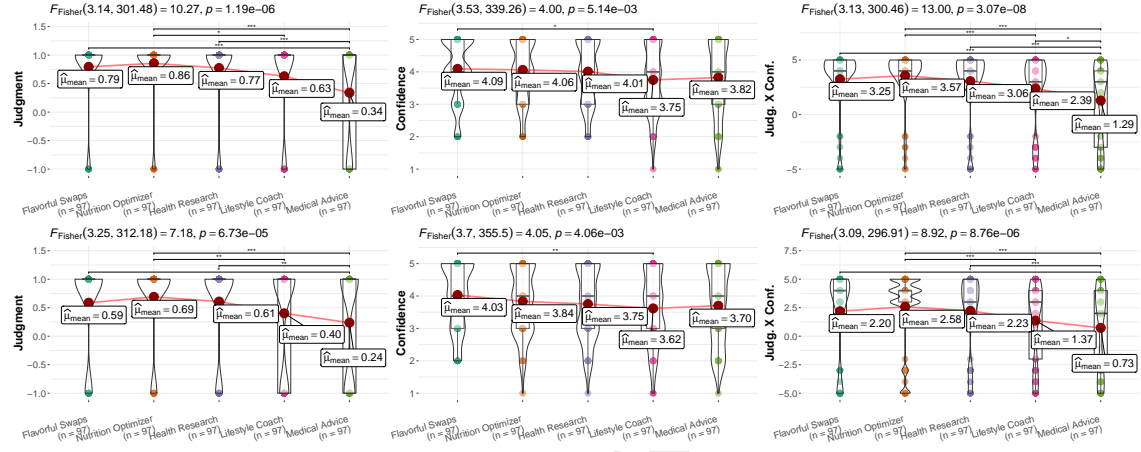


Fig. 4. Personal use case means and distributions of numerically converted Judgment, Confidence, and Judgment×Confidence. First row shows results for development decisions (Q1, Q2) and second row shows results for usage decisions (Q5, Q6). ANOVA results for use cases are shown above each panel. Within subject test was performed using Student's t-test with Holm correction. * denotes following significant p-values: $^{***}p < 0.001$; $^{**}p < 0.01$; $^{*}p < 0.05$. Use case names were shortened.

*Personal Use Cases.* In personal use scenarios, Digital Medical Advice AI ($M_{\text{DEV}} = 0.34, SD_{\text{DEV}} = 0.95$; $M_{\text{USAGE}} = 0.24, SD_{\text{USAGE}} = 0.98$), reflecting high risk level, consistently had lower acceptance across judgment types, compared to all other use cases. However, Customized Lifestyle Coach AI had the lowest confidence across both judgments ($M_{\text{DEV}} = 3.75, SD_{\text{DEV}} = 0.52$; $M_{\text{USAGE}} = 3.62, SD_{\text{USAGE}} = 0.78$). Nutrition Optimizer ($M_{\text{DEV}} = 0.86, SD_{\text{DEV}} = 0.92$; $M_{\text{USAGE}} = 0.69, SD_{\text{USAGE}} = 0.73$) had the highest mean acceptance across the two acceptability judgments. Interestingly, unlike the professional use cases, which had slightly higher acceptance for usage, personal use cases had lower acceptance for usage in general compared to development. Notably, many personal use cases had substantially high acceptance, especially compared to professional use cases.

*5.1.2 Use Case Variations.* When selecting use cases, we used two underlying variations: entry level of education required for professional use cases and EU AI risk levels for personal use cases. As risk levels and required education increased, we observe consistent negative effects on judgments, with personal use cases ($\beta_{\text{DEV}} = -0.11, p < .001$; $\beta_{\text{USAGE}} = -0.10, p < .001$) showing stronger effects compared to professional scenarios ($\beta_{\text{DEV}} = -0.08, p < .01$), where only development judgments were significantly associated. Again, acceptance for personal use cases (($intercept$)$_{\text{DEV}} = 1.02, p < .001$; ($intercept$)$_{\text{USAGE}} = 0.80, p < .001$) were higher than professional use cases ( ($intercept$)$_{\text{DEV}} = 0.43, p < .001$) as conveyed by the intercepts. Confidence ratings remained consistently high across all conditions (intercepts > 3.97), though they showed a small but significant decrease with increasing risk levels in personal use cases ($\beta_{\text{DEV}} = -0.08, p < .001$; $\beta_{\text{USAGE}} =$

| Use Case Variation | DEV($\beta(SE)$) | | USAGE($\beta(SE)$) | |
| --- | --- | --- | --- | --- |
| | Judgment | Confidence | Judgment | Confidence |
| **Professional** | | | | |
| Education Level | $-0.08^{**}$ (0.03) | $-0.00$ (0.02) | $0.00$ (0.03) | $-0.03$ (0.03) |
| Intercept | $0.43^{***}$ (0.10) | $3.97^{***}$ (0.10) | $0.17$ (0.10) | $4.04^{***}$ (0.10) |
| **Personal** | | | | |
| EU AI Risk Level | $-0.11^{***}$ (0.02) | $-0.08^{***}$ (0.02) | $-0.10^{***}$ (0.02) | $-0.09^{***}$ (0.02) |
| Intercept | $1.02^{***}$ (0.08) | $4.20^{***}$ (0.10) | $0.80^{***}$ (0.09) | $4.05^{***}$ (0.11) |

$^{***}p < 0.001; ^{**}p < 0.01; ^{*}p < 0.05$

Table 2. Results of `lmer` models analyzing the effects of use case variations on decisions, with subject-specific random effects. The data is separated into professional and personal categories. Use case variations are numerically coded from 1 (low risk/lowest entry level education) to 5 (high risk/highest entry level education). The table presents estimates, standard errors, and significance levels.

$-0.09, p < .001$), while professional use cases showed no significant impact on confidence. Thus, while education level required for entry does show some significant effect, the EU AI risk level had consistent effects with higher significance.

*5.1.3 Disagreements.* Observing the standard deviation and visualization of the distribution shows further understanding of possible disagreements among use cases. We compare the judgments weighted by confidence to understand not only the differences in judgment but also their strength. Interestingly, the use cases with four highest disagreements in both judgments were all professional uses in order of Telemarketer ($SD_{DEV} = 4.08$; $SD_{USAGE} = 4.21$), Elementary School Teacher ($SD_{DEV} = 3.99$; $SD_{USAGE} = 4.09$), Lawyer ($SD_{DEV} = 4.00$; $SD_{USAGE} = 3.99$), and Government Eligibility Interviewer AI ($SD_{DEV} = 3.96$; $SD_{USAGE} = 3.89$). These four use cases were followed by Digital Medical Advice AI ($SD_{DEV} = 3.80$, $SD_{USAGE} = 3.83$). The use cases with the lowest disagreements were surprisingly Nutrition Optimizer ($SD_{DEV} = 2.16$; $SD_{USAGE} = 3.03$) followed by IT Support Specialist AI ($SD_{DEV} = 3.08$; $SD_{USAGE} = 2.65$). These results underscore the general controversy of professional usages but also shows that acceptability is highly use case dependent.

## 5.2 RQ2. Impact of Personal Factors on Acceptability Judgment

*5.2.1 Demographic Factors.* In analyzing the demographic factors influencing judgments on the development and usage of AI use cases, several significant trends emerged from the data. First, we observed that the demographic variables had less than 0.5 correlation, except for age 65+ and Retired employment status. Across both categories of use cases, age was positively associated with confidence in acceptability judgment of usage for age 25-34 ($\beta = 0.41$, $p < 0.05$) and 55-64 ($\beta = 0.48$, $p < 0.05$). Race also had notable influences; specifically, Asian participants exhibited significantly lower confidence in both development and usage judgments ($\beta_{DEV} = -0.37$, $p < 0.01$; $\beta_{USAGE} = -0.33$, $p < 0.05$), particularly in professional contexts.

Gender emerged as a crucial determinant, with non-male participants consistently showing negative judgments across both development ($\beta = -0.29$, $p < 0.001$) and usage ($\beta = -0.33$, $p < 0.001$), indicating potential discrepancies in perception or experience with AI applications. Liberal views, especially among those identifying as strongly liberal, were associated with negative judgments across both categories of use cases ($\beta_{DEV} = -1.16$, $p < 0.05$; $\beta_{USAGE} = -1.51$, $p < 0.01$), suggesting a skeptical stance towards AI's prevalence and role. Employment hours also contributed, with individuals working 40+ hours per week displaying a positive association with development judgments ($\beta = 0.25$, $p < 0.05$), suggesting more exposure or reliance on AI use cases. High experience of discrimination chronicity was significantly related to lower acceptance of development, $\beta = -0.36$, $p < 0.05$. These findings highlight the significant role of demographic factors in shaping perceptions and attitudes toward AI technologies. See Table 23 in the Appendix for ANOVA results.

| Demographics | DEV ($\beta$ (SE)) | | | USAGE ($\beta$ (SE)) | | |
|---|---|---|---|---|---|---|
| | Judg. | Conf. | Judg.×Conf. | Judg. | Conf. | Judg.×Conf. |
| (Intercept) | **0.50*** (0.25) | **4.08***** (0.32) | 1.88 (1.09) | 0.51 (0.28) | **3.09***** (0.34) | 1.34 (1.23) |
| (Intercept)$_{Prof}$ | 0.24 (.38) | **4.59***** (.49) | 1.16 (1.66) | 0.71 (.41) | **3.74***** (.45) | 3.09 (1.74) |
| (Intercept)$_{Pers}$ | 0.66 (.30) | **3.76***** (.47) | 2.33 (1.34) | 0.25 (.43) | **2.61***** (.55) | −0.16 (1.92) |
| **Age** | | | | | | |
| 25-34 | −0.13 (0.13) | 0.04 (0.19) | −0.59 (0.58) | −0.22 (0.16) | **0.41*** (0.20) | −0.99 (0.69) |
| 55-64 | −0.03 (0.16) | 0.32 (0.23) | −0.14 (0.69) | 0.12 (0.19) | **0.48*** (0.24) | 0.40 (0.82) |
| 25-34$_{Pers}$ | −0.06 (.16) | 0.36 (.26) | −0.01 (.70) | −0.14 (.23) | **0.76*** (.30) | −0.10 (.103) |
| **Race** | | | | | | |
| Asian | 0.17 (0.10) | **−0.37**** (0.14) | 0.73 (0.44) | 0.10 (0.12) | **−0.33*** (0.15) | 0.42 (0.52) |
| Black | 0.07 (0.10) | 0.24 (0.14) | 0.49 (0.43) | 0.07 (0.12) | **0.32*** (0.15) | 0.53 (0.51) |
| Mixed | 0.19 (0.13) | 0.16 (0.18) | 0.85 (0.56) | −0.13 (0.15) | **0.43*** (0.19) | −0.12 (0.66) |
| Asian$_{Prof}$ | **0.40**** (.15) | **−0.41*** (.20) | **1.67*** (.65) | 0.20 (.16) | **−0.44*** (.19) | 0.59 (.68) |
| Asian$_{Pers}$ | 0.01 (.12) | **−0.54**** (.20) | −0.05 (.56) | 0.00 (.18) | −0.37 (.24) | 0.07 (.82) |
| Black$_{Pers}$ | 0.12 (.12) | 0.35 (.20) | 0.94 (.55) | 0.02 (.18) | **0.59*** (.23) | 0.65 (.80) |
| **Gender** | | | | | | |
| Non-male | **−0.29***** (0.07) | −0.05 (0.10) | **−1.29***** (0.32) | **−0.33***** (0.09) | 0.10 (0.11) | **−1.36***** (0.38) |
| Non-male$_{Prof}$ | **−0.48***** (.11) | 0.03 (.15) | **−2.11***** (.48) | **−0.52***** (.12) | 0.06 (.14) | **−2.25***** (.50) |
| **Political View** | | | | | | |
| Str. liberal | −0.19 (0.11) | −0.28 (0.16) | **−1.16*** (0.49) | **−0.34*** (0.13) | 0.14 (0.17) | **−1.51**** (0.59) |
| Str. Liberal$_{Prof}$ | −0.25 (.18) | 0.02 (.24) | −1.08 (.76) | **−0.42*** (.18) | **0.58*** (.22) | **−1.63*** (.79) |
| Str. Liberal$_{Pers}$ | −0.15 (.16) | **−0.53*** (.25) | −1.14 (.70) | −0.18 (.23) | −0.36 (.30) | −1.01 (.102) |
| Liberal$_{Pers}$ | −0.08 (.11) | **−0.52**** (.18) | −0.55 (.50) | −0.07 (.17) | −0.39 (.21) | −0.38 (.74) |
| **Education** | | | | | | |
| Advanced$_{Prof}$ | **0.43*** (.19) | −0.48 (.26) | 1.39 (0.83) | 0.31 (0.20) | −0.40 (0.24) | 1.03 (0.87) |
| **Employment** | | | | | | |
| 40+ hrs | **0.25*** (0.11) | 0.07 (0.16) | **1.13*** (0.51) | 0.09 (0.14) | 0.01 (0.17) | 0.73 (0.60) |
| **Discrimination** | | | | | | |
| High$_{Prof}$ | **−0.36*** (.18) | 0.14 (.25) | **−1.79*** (.79) | −0.13 (.19) | 0.27 (.23) | −0.39 (.82) |

*** $p < 0.001$; ** $p < 0.01$; * $p < 0.05$

Table 3. Coefficients, Standard Errors, and Significance of Demographic Factors. Models used decision metrics as dependent variables, demographic factors as independent variables, with random effects for subjects and use cases. Non-significant factors are excluded. The intercept represents the dominant demographic group (White, Christian, Male), the lowest natural ordering (age 18-24, Not Employed), and median scale values (Moderate, Associate's degree). Subscript $_{Prof}$ denote professional and $_{Pers}$ denote personal.

| AI Literacy | DEV ($\beta$ (SE)) | | | USAGE ($\beta$ (SE)) | | |
|---|---|---|---|---|---|---|
| | Judg. | Conf. | Judg.×Conf. | Judg. | Conf. | Judg.×Conf. |
| **Professional** | | | | | | |
| (Intercept) | 0.07 (0.30) | **3.13***** (0.33) | −0.49 (1.27) | −0.50 (0.32) | **3.46***** (0.34) | −2.44 (1.35) |
| AI Ethics | **−0.04*** (0.02) | 0.02 (0.02) | −0.15 (0.08) | −0.00 (0.02) | −0.01 (0.02) | −0.06 (0.08) |
| Gen AI Usage Freq. | **0.14**** (0.04) | −0.00 (0.05) | **0.59**** (0.19) | **0.18***** (0.05) | −0.07 (0.05) | **0.80***** (0.20) |
| Gen AI Limit. Familiarity | −0.09 (0.07) | **0.20*** (0.08) | −0.38 (0.28) | −0.06 (0.07) | **0.18*** (0.08) | −0.38 (0.30) |
| **Personal** | | | | | | |
| (Intercept) | **0.87***** (0.22) | **2.72***** (0.40) | **2.81**** (1.00) | 0.49 (0.30) | **2.44***** (0.45) | 1.17 (1.30) |
| AI Skills | 0.00 (0.01) | **0.06*** (0.02) | 0.07 (0.06) | 0.02 (0.02) | **0.06*** (0.03) | **0.15*** (0.08) |
| AI Ethics | **−0.05***** (0.01) | 0.01 (0.03) | **−0.23***** (0.07) | **−0.06**** (0.02) | **0.07*** (0.03) | **−0.26**** (0.09) |
| Gen AI Usage Freq. | **0.15***** (0.03) | 0.06 (0.06) | **0.62***** (0.15) | **0.19***** (0.05) | −0.01 (0.07) | **0.79***** (0.20) |
| Gen AI Limit. Familiarity | −0.06 (0.05) | 0.00 (0.09) | −0.25 (0.21) | −0.10 (0.06) | −0.10 (0.10) | −0.50 (0.28) |

*** $p < 0.001$; ** $p < 0.01$; * $p < 0.05$

Table 4. Coefficients, Standard Errors, and Significance of AI Literacy Factors. Models used decision metrics as dependent variables and AI literacy factors as independent variables, with random effects for subjects and use cases. Non-significant factors are excluded.

*5.2.2 AI Literacy.* We identified a correlation greater than 0.5 among three AI literacy aspects: awareness, usage, and evaluation. These were aggregated into a single factor, AI Skills. As shown in Table 4, understanding of AI Ethics was associated with lower acceptability for both personal ($\beta_{DEV} = -0.05, p < .001$; $\beta_{USAGE} = -0.23, p < .001$) and professional ($\beta_{DEV} = -0.04, p < .05$). However, across both categories of use cases, high Generative AI Usage Frequency

resulted in higher acceptance for both professional ($\beta_{\text{DEV}} = 0.14$, $p < .01$; $\beta_{\text{USAGE}} = 0.18$, $p < .001$) and usage acceptance ($\beta_{\text{DEV}} = 0.15$, $p < .001$; $\beta_{\text{USAGE}} = 0.19$, $p < .001$). Notably, for personal use cases, AI Skills was positively associated with confidence of judgments ($\beta = 0.06$, $p < .05$ for both development and usage) and judgment weighted by confidence ($\beta_{\text{USAGE}} = 0.15$, $p < .05$), while Generative AI Limitation Familiarity was positively associated with confidence for professional usage ($\beta_{\text{DEV}} = 0.20$, $p < .05$; $\beta_{\text{USAGE}} = 0.18$, $p < .05$ for usage). These results implicate that different understandings of and experiences with AI can impact judgments of acceptability, corroborating previous findings [43].

### 5.3 RQ3. Factors in Participant Rationale

To provide deeper insights to our analyses of acceptability judgments, we also examined open-text rationales (Q3) for judgments of development (Q1) and conditions for switching their decisions (Q4).
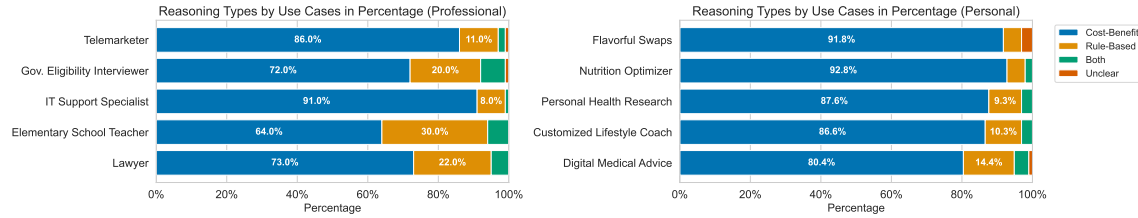


Fig. 5. Percentage of reasoning types (cost-benefit and rule-based) by use cases in participant provided rationales (Q3, Q4).

*5.3.1 Decision-making Types.* As we defined in § 4.2, we focus on two distinct reasoning types for decision-making, distinguished by consideration of outcome versus consideration of value inherent in action: cost-benefit reasoning (e.g., *"it gives more people access to medical advice and treatment"*, P365) and rule-based reasoning (e.g., should not be developed because *"human interaction is better"*, P249). As shown in Figure 5, generally, participants used more cost-benefit reasoning with highest percentage for IT Support Specialist (91.0%) and Nutrition Optimizer (92.8%). This result is particularly interesting as we observed these two use cases to have the lowest disagreement (see § 5.1.1) and signifies that perhaps unified reasoning type leads to less disagreement.

On the other hand, rationales for Elementary School Teacher AI contained most percentage of rule-based reasoning (30.0%) followed by Lawyer AI (22.0%). Interestingly, Elementary School Teacher AI and Lawyer AI were the use cases with the lowest development acceptability. These results suggest that less acceptable AI use cases might trigger more rule-based reasoning, which could, however, be due to the lack of consensus and rules around AI, especially those that are positive. this suggests the use of rule-based reasoning might correlate with judging AI use cases to be more unacceptable. These results suggest that the use of rule-based reasoning might correlate with judging AI use cases to be more unacceptable.
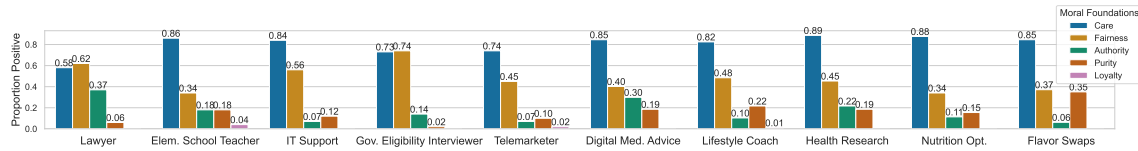


Fig. 6. Proportion of presence of moral foundations in participant's rationale responses (Q3, Q4) aggregated by use case.

*5.3.2 Moral Foundations.* Beyond reasoning types, we explored moral foundations to provide insights into what values are relevant for AI use case decisions (Figure 6). For example, P12 responded that Elementary School Teacher AI use case should be developed because it *"could give elementary schooling to children who are bed ridden..."*, which was annotated with both values of Care (focusing on the well-being and nurturing of bed-ridden children) and Fairness (focusing on fair access to education). Upon analysis, we observe that Care (i.e., dislike of pain of others, feelings of empathy and compassion toward others) was the most prevalent moral value in participants' rationales across the use cases. Of note, Care could impact the acceptability of AI use cases in both directions, as conveyed by P385 who noted that Customized Lifestyle Coach AI should be developed because "it may help improve some people's health" but would change their decision if "it caused harm to even one person." While our choice of medical domain in personal usage could have impacted the distribution, care was still the most prevalent considering only professional use cases.

Participants' rationales evoked fairness value more than care for two use cases, Lawyer AI (0.62) and Government Eligibility Interviewer (0.74) respectively. These results could be due to the characteristics of the use cases, such as their main purpose and function, as noted by P88, Government Eligibility Interviewer should be developed because *"it might be less biased and therefore more fair in it decisions (sic)"*. Authority was most apparent in participant rationales for Lawyer AI (0.37) and Purity for Flavorful Swaps (0.35).
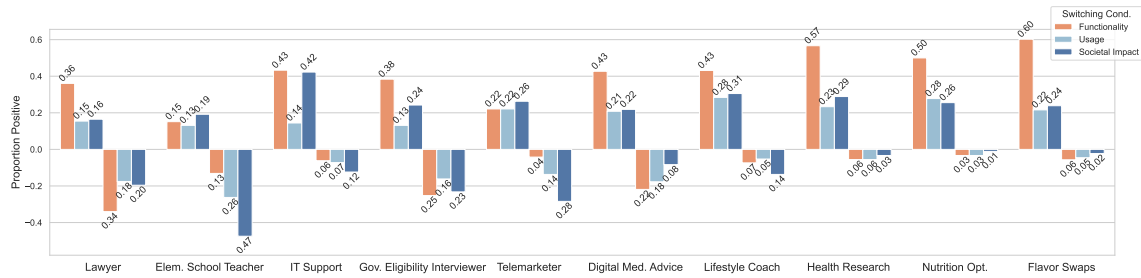


Fig. 7. Proportion of presence of switching conditions (Functionality, Usage, Societal Impact) mentioned in participant's switching conditions (Q4) aggregated by use case and divided into positive and negative development acceptability. The total positive proportion of the switching condition (indicated with same hue) by use case is the sum of both positive and negative bars.

*5.3.3 Switching Conditions.* We have thus far explored firmness of the decision through levels of confidence provided by users. We further explore the flexibility of participants' judgments to understand possible mitigation of disagreements through concerns expressed in conditions under which they would switch their decisions (Figure 7). Overall, Functionality (53.3%; e.g., Medical Adivce AI should not be developed if it "consistently or had a high percentage of failure to diagnose correctly.", P373) was the most commonly noted concern that would switch participants decision. This was followed by Societal Impact (42.3%; e.g., Lawyer AI should not be developed if "it puts too many human lawyers out of work", P97), Usage (31.8%; e.g., Government Eligibility Interviewer should be developed if "it was only used to read and screen applications but not for making decisions"), Not Applicable (0.02%; e.g., will not change decision).

Interestingly, Societal Impact (50.6%) was more prevalent followed by Functionality (46.4%) for professional use cases whereas Functionality (54.4%) was the most frequently mentioned concern for personal use cases. One of the reasons for frequent mention of Societal Impact in professional use cases could be due to concerns of labor replacement for professional use cases: as described by P296, if Elementary School Teacher AI *"was to replace teachers with the ai to save money"*, they would switch decision from should be developed to should not be developed. Moreover, we observe that

Societal Impact was the most prevalent concern for those (31%) who thought the use case was unacceptable across all use cases and was especially prevalent for Elementary School Teacher AI (0.47). These diverging results show the importance of understanding granular concerns by use cases to effectively address the relevant issues. Participants' responses that were identified as being both rule-based and concerned about societal impact expressed rules such as necessity of humanness as expressed by P22, "a human touch is 100% necessary" when discussing Elementary School AI, whether it be due to belief that "AI would lack empathy" (P5) or that "humans need human interactions" (P44).

| Rationale Factors | DEV ($\beta$ (SE)) | | | USAGE ($\beta$ (SE)) | | |
|---|---|---|---|---|---|---|
|  | Judg. | Conf. | Judg.×Conf. | Judg. | Conf. | Judg.×Conf. |
| (Intercept) | 0.07 (0.12) | 4.04*** (0.17) | 0.20 (0.46) | 0.15** (0.05) | 3.93*** (0.18) | −0.92 (0.68) |
| **Reasoning Type** | | | | | | |
| Cost-benefit | 0.33** (0.10) | 0.17 (0.15) | 1.35*** (0.39) | −0.10* (0.04) | 0.30 (0.16) | 2.07*** (0.52) |
| Rule-based | −0.47*** (0.09) | 0.44** (0.13) | −1.64*** (0.35) | −0.06 (0.04) | 0.26 (0.14) | −1.05* (0.47) |
| **Moral Value** | | | | | | |
| Care | 0.08 (0.05) | −0.21** (0.08) | 0.09 (0.20) | 0.02 (0.02) | −0.12 (0.08) | 0.44 (0.28) |
| Fairness | 0.10* (0.04) | −0.21** (0.07) | 0.47** (0.17) | 0.01 (0.02) | −0.18** (0.07) | 0.83*** (0.23) |
| Purity | 0.03 (0.06) | −0.12 (0.08) | 0.12 (0.21) | 0.01 (0.02) | −0.06 (0.09) | −0.13 (0.29) |
| Authority | −0.11 (0.06) | −0.18* (0.09) | −0.43 (0.22) | 0.01 (0.02) | −0.19* (0.09) | 0.17 (0.31) |
| **Switching Condition** | | | | | | |
| Functionality | −0.08* (0.04) | −0.05 (0.06) | −0.34* (0.15) | −0.01 (0.02) | −0.13* (0.06) | −0.05 (0.21) |
| Usage | 0.12*** (0.01) | 0.02* (0.01) | 0.56*** (0.02) | 0.24*** (0.00) | −0.04*** (0.01) | |
| Societal Impact | −0.11* (0.04) | 0.01 (0.06) | −0.49** (0.17) | 0.03 (0.02) | −0.13 (0.07) | −0.76*** (0.23) |

***$p < 0.001$; **$p < 0.01$; *$p < 0.05$

Table 5. Coefficients, Standard Errors, and Significance of Rationale Factors in Participants' Open-Text Responses. Models used decision metrics as dependent variables and rational factors and independent variables, with random effects for subjects and use cases. All rationale factors were coded as binary (presence: 0, 1). Non-significant factors are excluded.

*5.3.4 Rationale Factors and Judgment.* The results of analyzing the influence of rationale factors on judgments of acceptability indicated several significant effects (Table 5). We first verified that the factors had less than 0.5 correlation with each other. For acceptability of development, use of Cost-benefit reasoning in rationale was associated with increased acceptability ($\beta_{\mathsf{DEV}} = 0.33, p < .01$), while Rule-based reasoning was negatively associated with acceptability ($\beta_{\mathsf{DEV}} = -0.47, p < .001$). Interestingly, for usage, Cost-benefit reasoning was negatively associated with acceptability ($\beta_{\mathsf{USAGE}} = -0.10, p < .05$) but positively for judgment weighted by confidence ($\beta_{\mathsf{USAGE}} = 2.07, p < .001$). Rule-based reasoning indicated lower usage acceptability for judgment weighted by confidence, consistent with acceptability of development. In the moral value category, Fairness in rationales had positive associations with acceptance ($\beta_{\mathsf{DEV}} = 0.47, p < .01, \beta_{\mathsf{USAGE}} = 0.83, p < .001$; judgment weighted by confidence).

Analyzing switching conditions showed that participants who answered negatively to the development of use cases mentioned the concerns about Functionality ($\beta_{\mathsf{DEV}} = -0.08, p < .05$) and Societal Impact to switch their decisions ($\beta_{\mathsf{DEV}} = -0.11, p < .05$). However, those who were positive towards development of use case indicated emphasis on Usage ($\beta_{\mathsf{DEV}} = 0.12, p < .001$) as a condition to reverse their decisions. Thus, our findings highlight diverse perspectives on requirements and concerns of AI use cases, especially with varying perception of acceptability.

## 6 Conclusion and Discussion

We conducted a study to understand how and why laypeople perceive various AI use cases as acceptable or not. To achieve this, we developed a survey that gathered judgments and reasoning processes from 197 participants who were demographically diverse and had varying levels of experience with AI. Participants were asked to provide their judgments

on the acceptability of AI use cases, along with rationales for their decisions (e.g., "Should / Should not be developed, because...") and conditions that might change their decisions (e.g., "I would switch my decision if..."). The survey covered ten different AI use cases, spanning both personal and professional domains, and included varying levels of risk. Our findings revealed significant variation in the acceptability judgments and reasoning factors based on the domain, risk level, and participants' attributes, such as AI literacy and gender. We discuss the implications of these findings below.

*Use Case Perceptions and Disagreements.* In our study, we explored the varying acceptability of AI across different use cases. Generally, acceptance was lower in scenarios with higher educational requirements and greater EU AI risk levels. Professional use cases displayed more variability, notably with Elementary School Teacher AI, which was uniquely unacceptable. This underscores the necessity for further research into how AI should be developed and integrated, as well as what skills it should have, particularly in fields where empathy and care are crucial [15, 38, 82]. In addition, prior FAccT research have also highlighted how AI practitioners desire understanding lay people's perception on AI fairness in specific use cases [25, 26, 65, 69]. Drawing from prior HCI and AI research [19, 27, 45], future FAccT researchers and practitioners should explore how to meaningfully connect lay people's use case perceptions with AI developers' workflows.

While prior research has emphasized understanding AI consequences [41] and providing tools and processes to uncover impact [16, 24, 78], our findings reveal a greater presence of rule-based reasoning in contentious use cases, suggesting a need for diverse approaches to understanding AI beyond mere consequence anticipation. Moreover, while care was generally predominant, we observed that fairness gained prominence in Lawyer AI and Government Eligibility Interviewer AI. This variability underscores the importance of considering values in AI evaluation and training [9, 12], rather than solely emphasizing functionality, which is the current trend in AI research [13]. Additionally, societal impact considerations were more evident in unacceptable use cases, emphasizing the necessity for implementing safety guardrails when deploying AI with significant social implications [67].

*Demographics and AI Literacy.* In line with prior work [42, 53], our results highlighted significant differences among demographic groups and perceived acceptance of use cases, especially for professional use (§5.2). Non-majority demographic groups, especially non-male gender groups, found both personal and professional use cases less acceptable. Those experiencing high discrimination chronicity also found professional use less acceptable. Our findings provide future FAccT research with valuable empirical insights on AI integration in contexts such as workplace, where marginalized worker's agency, earning, and occupational well-being are disproportionately affected [5, 52].

Furthermore, our work highlighted a potential polarization on perceptions of AI among workers as those with 40+ hours employment and those who had advanced degrees were more positive towards AI use cases, suggesting that the relationship stakeholders have to AI and jobs might influence acceptability. Such a concern was expressed by one of our participants who opposed development of Telemarketer AI because it *"overlaps with my industry, and hence serves as a threat to my job security"* (P35). Thus, our results corroborate the need to further explore methods to include diverse workers and various stakeholders into the discussion of workplace AI integration and development [20, 30]. We also found that frequent AI usage increased acceptance, while understanding AI ethics and limitations decreased acceptance. This suggests that balanced AI awareness and education, encompassing usage, skills, and ethics, could guide and improve decision-making [61], e.g., through educational interventions targeting AI skills and ethical implication literacy (e.g., [64, 80]).

*Rationales.* Through analyzing participants' rationales, we observed an interesting pattern with higher cost-benefit reasoning use cases with least disagreement and more prevalence of rule-based reasoning for those with higher disagreement (§5.3). These results suggest different valuation systems employed by participants perhaps leading to different

conclusions and suggests some use cases have beyond simple utilitarian implication for society, which should be explored more carefully in future studies. Building upon our empirical findings, future work could develop tools and interventions which encourage specific types of acceptability reasoning such as rule and value based [70] or cost-benefit analyses [46].

However, as our study was limited to the two reasoning type categories, expanding this analysis would be essential for future work including finding ways to classify what features people are considering in their decisions, how the weights on those features impact what kind of decision-making strategy they will use, and whether there are other ways to understand their decision strategies beyond our current classification. Future FAccT research can build upon these further understandings to guide policy making and consensus building. For example, in addition to conducting surveys with single participant, future work can explore how group discussions and deliberations shape communities' collective understanding of AI impacts (e.g., [28, 32, 44, 45, 84])

## References

[1] 2023. Executive Order on the Safe, Secure, and Trustworthy Development and Use of Artificial Intelligence. https://www.whitehouse.gov/briefing-room/presidential-actions/2023/10/30/executive-order-on-the-safe-secure-and-trustworthy-development-and-use-of-artificial-intelligence/. Accessed: 2024-01-02.

[2] 2023. How Do People Feel About AI? A Nationally Representative Survey of Public Attitudes to Artificial Intelligence in Britain.

[3] ACL, Ethic Policy. 2023. ACL 2023 Responsible NLP Research and Ethics Policy. (March 2023). https://2023.aclweb.org/calls/main_conference/#ethics-policy

[4] Ada Lovelace Institute. 2022. Looking before we leap: Expanding ethical review processes for AI and data science research. (December 2022). https://www.adalovelaceinstitute.org/wp-content/uploads/2022/12/Ada-Lovelace-Institute-Looking-before-we-leap-Dec-2022.pdf

[5] Carlos-Maria Alcover, Dina Guglielmi, Marco Depolo, and Greta Mazzetti. 2021. "Aging-and-Tech Job Vulnerability": A proposed framework on the dual impact of aging and AI, robotics, and automation among older workers. *Organizational Psychology Review* 11, 2 (2021), 175–201.

[6] Carolyn Ashurst, Solon Barocas, Rosie Campbell, Deborah Raji, and Stuart Russell. 2020. Navigating the Broader Impacts of AI Research. NeurIPS workshop. https://ai-broader-impacts-workshop.github.io/

[7] Mohammad Atari, Jonathan Haidt, Jesse Graham, Sena Koleva, Sean T Stevens, and Morteza Dehghani. 2023. Morality beyond the WEIRD: How the nomological network of morality varies across cultures. *Journal of Personality and Social Psychology* (2023).

[8] Edmond Awad, Sohan Dsouza, Richard Kim, Jonathan Schulz, Joseph Henrich, Azim Shariff, Jean-François Bonnefon, and Iyad Rahwan. 2018. The moral machine experiment. *Nature* 563, 7729 (2018), 59–64.

[9] Solon Barocas, Anhong Guo, Ece Kamar, Jacquelyn Krones, Meredith Ringel Morris, Jennifer Wortman Vaughan, W. Duncan Wadsworth, and Hanna Wallach. 2021. Designing Disaggregated Evaluations of AI Systems: Choices, Considerations, and Tradeoffs. In *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society* (Virtual Event, USA) *(AIES '21)*. Association for Computing Machinery, New York, NY, USA, 368–378. https://doi.org/10.1145/3461702.3462610

[10] Solon Barocas and Andrew D Selbst. 2016. Big data's disparate impact. *Calif. L. Rev.* 104 (2016), 671.

[11] Michael S Bernstein, Margaret Levi, David Magnus, Betsy A Rajala, Debra Satz, and Quinn Waeiss. 2021. Ethics and society review: Ethics reflection as a precondition to research funding. *Proceedings of the National Academy of Sciences* 118, 52 (2021), e2117261118.

[12] Eshta Bhardwaj, Harshit Gujral, Siyi Wu, Ciara Zogheib, Tegan Maharaj, and Christoph Becker. 2024. Machine learning data practices through a data curation lens: An evaluation framework. In *The 2024 ACM Conference on Fairness, Accountability, and Transparency*. 1055–1067.

[13] Abeba Birhane, Pratyusha Kalluri, Dallas Card, William Agnew, Ravit Dotan, and Michelle Bao. 2022. The values encoded in machine learning research. In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*. 173–184.

[14] Su Lin Blodgett, Solon Barocas, Hal Daumé III, and Hanna Wallach. 2020. Language (technology) is power: A critical survey of" bias" in nlp. *arXiv preprint arXiv:2005.14050* (2020).

[15] Jana Schaich Borg and Hannah Read. 2024. What Is Required for Empathic AI? It Depends, and Why That Matters for AI Developers and Users. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*, Vol. 7. 1306–1318.

[16] Zana Buçinca, Chau Minh Pham, Maurice Jakesch, Marco Tulio Ribeiro, Alexandra Olteanu, and Saleema Amershi. 2023. AHA!: Facilitating AI Impact Assessment by Generating Examples of Harms. *arXiv preprint arXiv:2306.03280* (2023).

[17] Center for Advanced Study in the Behavioral Sciences. 2020. Ethics & Society Review · Stanford University. (2020). https://casbs.stanford.edu/ethics-society-review-stanford-university

[18] Richard J Chen, Judy J Wang, Drew FK Williamson, Tiffany Y Chen, Jana Lipkova, Ming Y Lu, Sharifa Sahai, and Faisal Mahmood. 2023. Algorithmic fairness in artificial intelligence for medicine and healthcare. *Nature biomedical engineering* 7, 6 (2023), 719–742.

[19] Hao-Fei Cheng, Ruotong Wang, Zheng Zhang, Fiona O'connell, Terrance Gray, F Maxwell Harper, and Haiyi Zhu. 2019. Explaining decision-making algorithms through UI: Strategies to help non-expert stakeholders. In *Proceedings of the 2019 chi conference on human factors in computing systems.*

1–12.

[20] EunJeong Cheon. 2023. Powerful Futures: How a Big Tech Company Envisions Humans and Technologies in the Workplace of the Future. *Proc. ACM Hum.-Comput. Interact.* 7, CSCW2, Article 312 (Oct. 2023), 35 pages. https://doi.org/10.1145/3610103

[21] Vanessa Cheung, Maximilian Maier, and Falk Lieder. 2024. Measuring the decision process in (moral) dilemmas: Self-report measures of reliance on rules, cost-benefit reasoning, intuition, & deliberation.

[22] Fiery Cushman. 2013. Action, outcome, and value: A dual-system framework for morality. *Personality and social psychology review* 17, 3 (2013), 273–292.

[23] CVPR, Ethics Guidelines. 2023. CVPR 2024 Ethics Guidelines for Authors. (November 2023). https://cvpr2023.thecvf.com/Conferences/2023/EthicsGuidelines

[24] Wesley Hanwen Deng, Solon Barocas, and Jennifer Wortman Vaughan. 2024. Supporting Industry Computing Researchers in Assessing, Articulating, and Addressing the Potential Negative Societal Impact of Their Work. *arXiv preprint arXiv:2408.01057* (2024).

[25] Wesley Hanwen Deng, Boyuan Guo, Alicia Devrio, Hong Shen, Motahhare Eslami, and Kenneth Holstein. 2023. Understanding Practices, Challenges, and Opportunities for User-Engaged Algorithm Auditing in Industry Practice. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*. 1–18.

[26] Wesley Hanwen Deng, Manish Nagireddy, Michelle Seng Ah Lee, Jatinder Singh, Zhiwei Steven Wu, Kenneth Holstein, and Haiyi Zhu. 2022. Exploring how machine learning practitioners (try to) use fairness toolkits. In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*. 473–484.

[27] Wesley Hanwen Deng, Claire Wang, Howard Ziyu Han, Jason I Hong, Kenneth Holstein, and Motahhare Eslami. 2025. WeAudit: Scaffolding User Auditors and AI Practitioners in Auditing Generative AI. *arXiv preprint arXiv:2501.01397* (2025).

[28] Alicia DeVos, Aditi Dhabalia, Hong Shen, Kenneth Holstein, and Motahhare Eslami. 2022. Toward User-Driven Algorithm Auditing: Investigating users' strategies for uncovering harmful algorithmic behavior. In *Proceedings of the 2022 CHI conference on human factors in computing systems*. 1–19.

[29] Julia Evetts. 2006. Introduction: Trust and professionalism: Challenges and occupational changes. , 515–531 pages.

[30] Sarah E. Fox, Vera Khovanskaya, Clara Crivellaro, Niloufar Salehi, Lynn Dombrowski, Chinmay Kulkarni, Lilly Irani, and Jodi Forlizzi. 2020. Worker-Centered Design: Expanding HCI Methods for Supporting Labor. 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, 1–8. https://doi.org/10.1145/3334480.3375157

[31] Delaram Golpayegani, Harshvardhan J. Pandit, and Dave Lewis. 2023. To Be High-Risk, or Not To Be—Semantic Specifications and Implications of the AI Act's High-Risk AI Applications and Harmonised Standards. In *Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency* (Chicago, IL, USA) *(FAccT '23)*. Association for Computing Machinery, New York, NY, USA, 905–915. https://doi.org/10.1145/3593013.3594050

[32] Mitchell L Gordon, Michelle S Lam, Joon Sung Park, Kayur Patel, Jeff Hancock, Tatsunori Hashimoto, and Michael S Bernstein. 2022. Jury learning: Integrating dissenting voices into machine learning models. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*. 1–19.

[33] Jesse Graham, Brian A Nosek, Jonathan Haidt, Ravi Iyer, Spassena Koleva, and Peter H Ditto. 2011. Mapping the moral domain. *Journal of personality and social psychology* 101, 2 (2011), 366.

[34] Jesse Graham, Brian A Nosek, Jonathan Haidt, Ravi Iyer, Koleva Spassena, and Peter H Ditto. 2008. Moral foundations questionnaire. *Journal of Personality and Social Psychology* (2008).

[35] Brent Hecht, Lauren Wilcox, Jeffrey P Bigham, Johannes Schöning, Ehsan Hoque, Jason Ernst, Yonatan Bisk, Luigi De Russis, Lana Yarosh, Bushra Anjum, et al. 2021. It's time to do something: Mitigating the negative impacts of computing through a change to the peer review process. *arXiv preprint arXiv:2112.09544* (2021).

[36] Viviane Herdel, Sanja Šćepanović, Edyta Bogucka, and Daniele Quercia. 2024. ExploreGen: Large language models for envisioning the uses and risks of AI technologies. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*, Vol. 7. 584–596.

[37] ICML, Publication Ethics. 2023. International Conference on Machine Learning, Publication Ethics. (November 2023). https://icml.cc/Conferences/2023/PublicationEthics

[38] Anna Kawakami, Jordan Taylor, Sarah Fox, Haiyi Zhu, and Ken Holstein. 2024. AI Failure Loops in Feminized Labor: Understanding the Interplay of Workplace AI and Occupational Devaluation. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*, Vol. 7. 683–683.

[39] Selin Kesebir and Jonathan Haidt. 2010. Morality (in Handbook of social psychology). *HANDBOOK OF SOCIAL PSYCHOLOGY, 5th Ed., S. Fiske, D. Gilbert, & G. Lindzey, eds., Forthcoming* (2010).

[40] Kimon Kieslich, Nicholas Diakopoulos, and Natali Helberger. 2023. Anticipating Impacts: Using Large-Scale Scenario Writing to Explore Diverse Implications of Generative AI in the News Environment. *arXiv preprint arXiv:2310.06361* (2023).

[41] Kimon Kieslich, Natali Helberger, and Nicholas Diakopoulos. 2024. My Future with My Chatbot: A Scenario-Driven, User-Centric Approach to Anticipating AI Impacts. In *Proceedings of the 2024 ACM Conference on Fairness, Accountability, and Transparency* (Rio de Janeiro, Brazil) *(FAccT '24)*. Association for Computing Machinery, New York, NY, USA, 2071–2085. https://doi.org/10.1145/3630106.3659026

[42] Sara Kingsley, Jiayin Zhi, Wesley Hanwen Deng, Jaimie Lee, Sizhe Zhang, Motahhare Eslami, Kenneth Holstein, Jason I Hong, Tianshi Li, and Hong Shen. 2024. Investigating What Factors Influence Users' Rating of Harmful Algorithmic Bias and Discrimination. In *Proceedings of the AAAI Conference on Human Computation and Crowdsourcing*, Vol. 12. 75–85.

[43] Max F. Kramer, Jana Schaich Borg, Vincent Conitzer, and Walter Sinnott-Armstrong. 2018. When Do People Want AI to Make Decisions?. In *Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society* (New Orleans, LA, USA) *(AIES '18)*. Association for Computing Machinery, New York, NY, USA, 204–209. https://doi.org/10.1145/3278721.3278752

[44] Tzu-Sheng Kuo, Quan Ze Chen, Amy X Zhang, Jane Hsieh, Haiyi Zhu, and Kenneth Holstein. 2024. PolicyCraft: Supporting Collaborative and Participatory Policy Design through Case-Grounded Deliberation. *arXiv preprint arXiv:2409.15644* (2024).

[45] Min Kyung Lee, Daniel Kusbit, Anson Kahng, Ji Tae Kim, Xinran Yuan, Allissa Chan, Daniel See, Ritesh Noothigattu, Siheon Lee, Alexandros Psomas, and Ariel D. Procaccia. 2019. WeBuildAI: Participatory Framework for Algorithmic Governance. *Proc. ACM Hum.-Comput. Interact.* 3, CSCW, Article 181 (Nov. 2019), 35 pages. https://doi.org/10.1145/3359283

[46] Jing-Jing Li, Valentina Pyatkin, Max Kleiman-Weiner, Liwei Jiang, Nouha Dziri, Anne G. E. Collins, Jana Schaich Borg, Maarten Sap, Yejin Choi, and Sydney Levine. 2024. SafetyAnalyst: Interpretable, transparent, and steerable LLM safety moderation. *arXiv* (2024). http://arxiv.org/abs/2410.16665

[47] Hsuan-Tien Lin, Maria-Florina Balcan, Raia Hadsell, and Marc'Aurelio Ranzato. 2020. Getting Started with NeurIPS 2020. NeurIPS blog. https://neuripsconf.medium.com/getting-started-with-neurips-2020-e350f9b39c28

[48] Jacob Metcalf, Emanuel Moss, Elizabeth Anne Watkins, Ranjit Singh, and Madeleine Clare Elish. 2021. Algorithmic impact assessments and accountability: The co-construction of impacts. In *Proceedings of the 2021 ACM conference on fairness, accountability, and transparency*. 735–746.

[49] Eli Michaels, Marilyn Thomas, Alexis Reeves, Melisa Price, Rebecca Hasson, David Chae, and Amani Allen. 2019. Coding the Everyday Discrimination Scale: implications for exposure assessment and associations with hypertension and depression among a cross section of mid-life African American women. *J Epidemiol Community Health* 73, 6 (2019), 577–584.

[50] Microsoft. 2022. Microsoft Responsible AI Impact Assessment Guide. (June 2022). https://blogs.microsoft.com/wp-content/uploads/prod/sites/5/2022/06/Microsoft-RAI-Impact-Assessment-Guide.pdf

[51] Microsoft. 2022. Microsoft Responsible AI Impact Assessment Template. (June 2022). https://blogs.microsoft.com/wp-content/uploads/prod/sites/5/2022/06/Microsoft-RAI-Impact-Assessment-Template.pdf

[52] Joy Ming, Lucy Pei, Rama Adithya Varanasi, Anna Kawakami, Nervo Verdezoto, and EunJeong Cheon. 2024. Labor, Visibility, and Technology: Weaving Together Academic Insights and On-Ground Realities. In *Companion Publication of the 2024 Conference on Computer-Supported Cooperative Work and Social Computing* (San Jose, Costa Rica) *(CSCW Companion '24)*. Association for Computing Machinery, New York, NY, USA, 708–711. https://doi.org/10.1145/3678884.3681827

[53] Jimin Mun, Liwei Jiang, Jenny Liang, Inyoung Cheong, Nicole DeCario, Yejin Choi, Tadayoshi Kohno, and Maarten Sap. 2024. Particip-AI: A Democratic Surveying Framework for Anticipating Future AI Use Cases, Harms and Benefits. arXiv:2403.14791 [cs.CY] https://arxiv.org/abs/2403.14791

[54] National Artificial Intelligence Research Resource Task Force. 2023. Strengthening and Democratizing the U.S. Artificial Intelligence Innovation Ecosystem. (Janurary 2023). https://www.ai.gov/wp-content/uploads/2023/01/NAIRR-TF-Final-Report-2023.pdf

[55] National Institute of Standards and Technology. 2023. *Artificial Intelligence Risk Management Framework (AI RMF 1.0)*. https://nvlpubs.nist.gov/nistpubs/ai/NIST.AI.100-1.pdf

[56] Alexandra Olteanu, Michael Ekstrand, Carlos Castillo, and Jina Suh. 2023. Responsible AI Research Needs Impact Statements Too. *arXiv preprint arXiv:2311.11776* (2023).

[57] Partnership on AI. 2021. Managing the Risks of AI Research: Six Recommendations for Responsible Publication. (2021). https://partnershiponai.org/paper/responsible-publication-recommendations/

[58] OpenAI. 2022. OpenAI: Our approach to alignment research. https://openai.com/blog/our-approach-to-alignment-research/

[59] European Parliament. 2023. Artificial Intelligence Act: deal on comprehensive rules for trustworthy AI. https://www.europarl.europa.eu/news/en/press-room/20231206IPR15699/artificial-intelligence-act-deal-on-comprehensive-rules-for-trustworthy-ai

[60] Giada Pistilli, Carlos Muñoz Ferrandis, Yacine Jernite, and Margaret Mitchell. 2023. Stronger together: on the articulation of ethical charters, legal tools, and technical documentation in ML. In *Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency*. 343–354.

[61] Inioluwa Deborah Raji, Morgan Klaus Scheuerman, and Razvan Amironesei. 2021. You Can't Sit With Us: Exclusionary Pedagogy in AI Ethics Education. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency* (Virtual Event, Canada) *(FAccT '21)*. Association for Computing Machinery, New York, NY, USA, 515–525. https://doi.org/10.1145/3442188.3445914

[62] Dillon Reisman, Jason Schultz, Kate Crawford, and Meredith Whittaker. 2018. Algorithmic Impact Assessments: A Practical Framework for Public Agency. *AI Now* (2018).

[63] Maarten Sap, Dallas Card, Saadia Gabriel, Yejin Choi, and Noah A Smith. 2019. The risk of racial bias in hate speech detection. In *Proceedings of the 57th annual meeting of the association for computational linguistics*. 1668–1678.

[64] Hong Shen, Wesley H Deng, Aditi Chattopadhyay, Zhiwei Steven Wu, Xu Wang, and Haiyi Zhu. 2021. Value cards: An educational toolkit for teaching social impacts of machine learning through deliberation. In *Proceedings of the 2021 ACM conference on fairness, accountability, and transparency*. 850–861.

[65] Jessie J Smith, Lex Beattie, and Henriette Cramer. 2023. Scoping fairness objectives and identifying fairness metrics for recommender systems: The practitioners' perspective. In *Proceedings of the ACM Web Conference 2023*. 3648–3659.

[66] Data & Society. 2023. Data & Society Announces the Launch of its Algorithmic Impact Methods Lab. https://datasociety.net/algorithmic-impact-methods-lab/

[67] Irene Solaiman. 2023. The Gradient of Generative AI Release: Methods and Considerations. In *Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency* (Chicago, IL, USA) *(FAccT '23)*. Association for Computing Machinery, New York, NY, USA, 111–122. https://doi.org/10.1145/3593013.3593981

[68] Irene Solaiman, Zeerak Talat, William Agnew, Lama Ahmad, Dylan Baker, Su Lin Blodgett, Canyu Chen, Hal Daumé III, Jesse Dodge, Isabella Duan, et al. 2023. Evaluating the social impact of generative ai systems in systems and society. *arXiv preprint arXiv:2306.05949* (2023).

[69] Nasim Sonboli, Jessie J Smith, Florencia Cabral Berenfus, Robin Burke, and Casey Fiesler. 2021. Fairness and transparency in recommendation: The users' perspective. In *Proceedings of the 29th ACM Conference on User Modeling, Adaptation and Personalization*. 274–279.

[70] Taylor Sorensen, Liwei Jiang, Jena Hwang, Sydney Levine, Valentina Pyatkin, Peter West, Nouha Dziri, Ximing Lu, Kavel Rao, Chandra Bhagavatula, Maarten Sap, John Tasioulas, and Yejin Choi. 2024. Value Kaleidoscope: Engaging AI with Pluralistic Human Values, Rights, and Duties. In *AAAI*. https://arxiv.org/abs/2309.00779

[71] Harini Suresh, Emily Tseng, Meg Young, Mary Gray, Emma Pierson, and Karen Levy. 2024. Participation in the age of foundation models. In *The 2024 ACM Conference on Fairness, Accountability, and Transparency*. 1609–1621.

[72] Lennart G Svensson. 2006. Professional occupations and status: A sociological study of professional occupations, status and trust. (2006).

[73] The Ada Lovelace Insitute. 2022. Algorithmic impact assessment: a case study in healthcare. (Feb 2022). https://www.adalovelaceinstitute.org/report/algorithmic-impact-assessment-case-study-healthcare/

[74] The White House. 2023. National Artificial Intelligence Research Resource Task Force Releases Final Report. (Janurary 2023). https://www.whitehouse.gov/ostp/news-updates/2023/01/24/national-artificial-intelligence-research-resource-task-force-releases-final-report/

[75] The New York Times. 2024. Will Chatbots Teach Your Children? (2024). https://www.nytimes.com/2024/01/11/technology/ai-chatbots-khan-education-tutoring.html

[76] Kent Walker and Marian Croak. 2021. An update on our progress in responsible AI innovation. (2021). https://blog.google/technology/ai/update-our-progress-responsible-ai-innovation/

[77] Bingcheng Wang, Pei-Luen Patrick Rau, and Tianyi Yuan. 2023. Measuring user competence in using artificial intelligence: validity and reliability of artificial intelligence literacy scale. *Behaviour & information technology* 42, 9 (2023), 1324–1337.

[78] Zijie J. Wang, Chinmay Kulkarni, Lauren Wilcox, Michael Terry, and Michael Madaio. 2024. Farsight: Fostering Responsible AI Awareness During AI Application Prototyping. In *Proceedings of the CHI Conference on Human Factors in Computing Systems*. 1–40.

[79] Laura Weidinger, Jonathan Uesato, Maribeth Rauh, Conor Griffin, Po-Sen Huang, John Mellor, Amelia Glaese, Myra Cheng, Borja Balle, Atoosa Kasirzadeh, et al. 2022. Taxonomy of risks posed by language models. In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*. 214–229.

[80] Richmond Y Wong and Tonya Nguyen. 2021. Timelines: A world-building activity for values advocacy. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*. 1–15.

[81] Nahathai Wongpakaran, Tinakon Wongpakaran, Danny Wedding, and Kilem L Gwet. 2013. A comparison of Cohen's Kappa and Gwet's AC1 when calculating inter-rater reliability coefficients: a study conducted with personality disorder samples. *BMC medical research methodology* 13 (2013), 1–7.

[82] Yiying Wu, Jung-Joo Lee, Ajit G. Pillai, Janghee Cho, Naseem Ahmadpour, Virpi Roto, Thida Sachathep, Jiashuo Liu, Mouna Sawan, Dongjin Song, Martina Čaić, Lucas Cheng, Renxuan Liu, Sarah Kettley, Luis Soares, Kazjon Grace, and Thomas Astell-Burt. 2024. Collective Imaginaries for the Futures of Care Work. In *Companion Publication of the 2024 Conference on Computer-Supported Cooperative Work and Social Computing* (San Jose, Costa Rica) *(CSCW Companion '24)*. Association for Computing Machinery, New York, NY, USA, 732–735. https://doi.org/10.1145/3678884.3681838

[83] Chunpeng Zhai, Santoso Wibowo, and Lily D Li. 2024. The effects of over-reliance on AI dialogue systems on students' cognitive abilities: a systematic review. *Smart Learning Environments* 11, 1 (2024), 28.

[84] Angie Zhang, Olympia Walker, Kaci Nguyen, Jiajun Dai, Anqing Chen, and Min Kyung Lee. 2023. Deliberating with AI: Improving Decision-Making for the Future through Participatory AI Design and Stakeholder Deliberation. *Proc. ACM Hum.-Comput. Interact.* 7, CSCW1, Article 125 (April 2023), 32 pages. https://doi.org/10.1145/3579601

[85] Dora Zhao, Angelina Wang, and Olga Russakovsky. 2021. Understanding and evaluating racial biases in image captioning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 14830–14840.

## A  Additional Survey Details

### A.1  Survey Questions

In our survey we ask participants to read one or more descriptions of AI use cases and to make two judgments: 1) "Do you think a technology like this should exist?" (Q1) and 2) "If the *<use case>* exists, would you use its services?" (Q5). A question to indicate their level of confidence is asked following each question (Q5, Q6). The participants are asked to both elaborate on their decisions (Q3) and specify the conditions under which they would switch their decisions (Q4). The detailed wordings for the questions are shown in Table 6.

Following the main use case questions in both the main and second study, we also asked participants questions about their demographics and literacy levels in AI, and the questions can be found in Table 7 and 8 respectively.

Lastly, while not included in the main text, we asked participants 3 questionnaires about decision making styles to explore the relationship between several decision making styles and the actual decisions of the participants. These

included: (1) Oxford Utilitarianism Scale, (2) Toronto Empathy Questionnaire and (3) Moral Foundations Questionnaire. The decision making style questions can be found in Table 9, 10 and 11 respectively.

| Question ID | Question | Answer Type |
|---|---|---|
| AI Perception Question (Before) | | |
| AI Perception Before | Overall, how does the growing presence of artificial intelligence (AI) in daily life and society make you feel? | 5 Point Likert Scale |
| Part 1 Specific Questions | | |
| UCX-1 | Do you think a technology like this should be developed? | Yes/No |
| UCX-2 | How confident are you in your above answer? | 5 Point Likert Scale |
| UCX-3Y | Please complete the following: [Use Case] should be developed because... | Text |
| UCX-3N | Please complete the following: [Use Case] should not be developed because... | Text |
| UCX-4Y | Under what circumstances would you switch your decision from [UCX-2 Answer] should be developed to should not be developed? | Text |
| UCX-4N | Under what circumstances would you switch your decision from [UCX-2 Answer] should not be developed to should be developed? | Text |
| UCX-5 | If [Use Case] exists, would you ever use its services (answer yes, even if you think you would use it very infrequently)? | Yes/No |
| UCX-6 | How confident are you in your above answer? | 5 Point Likert Scale |
| AI Perception Question (After) | | |
| AI Perception After | Before we continue, we'd like to get your thoughts on AI one more time. Overall, how does the growing presence of artificial intelligence (AI) in daily life and society make you feel? | 5 Point Likert Scale |

Table 6. Main Study Specific Question, the "X" in Question IDs is a placeholder for the use case number, which ranges from 1 to 5, for the 5 use cases in the jobs and personal use cases respectively.

| Question ID | Question |
|---|---|
| D-Q1 | How old are you? |
| D-Q2 | Choose one or more races that you consider yourself to be |
| D-Q3 | Do you identify as transgender? |
| D-Q4 | How would you describe your gender identity? |
| D-Q5 | How would you describe your sexual orientation? |
| D-Q6 | What is your present religion or religiosity, if any? |
| D-Q7 | In general, would you describe your political views as… |
| D-Q8 | What is the highest level of education you have completed? |
| D-Q9 | In which country have you lived in the longest? |
| D-Q10 | What other countries have you lived in for at least 6 months? |
| D-Q11 | Which of the following categories best describe your employment status? |
| D-Q12 | How would you describe the industry your job would be in? (Select all that apply) |
| D-Q13 | Do you identify with any minority, disadvantaged, demographic, or other specific groups? If so, which one(s)? (E.g., racial, gender identity, sexuality, disability, immigrant, veteran, etc.); use commas to separate groups. |
| D-Q14 | (Optional) What are some things that you are most concerned about lately? |
| $D-Q15_1$ | In your day-to-day life how often have any of the following things happened to you? You are treated with less courtesy or respect than other people |
| $D-Q15_2$ | In your day-to-day life how often have any of the following things happened to you? You receive poorer service than other people at restaurants or stores |
| $D-Q15_3$ | In your day-to-day life how often have any of the following things happened to you? People act as if they think you are not smart |
| $D-Q15_4$ | In your day-to-day life how often have any of the following things happened to you? People act as if they are afraid of you |
| $D-Q15_5$ | In your day-to-day life how often have any of the following things happened to you? You are threatened or harassed |
| D-Q16 | If the answer to Q15 is "A few times a year" or more frequently to at least one of the statements, what do you think is the main reason for these experiences? (Select all that apply) |

Table 7. Demographic Questions

| Question ID | Question |
|---|---|
| AI-Q1 | I can identify the AI technology employed in the applications and products I use. |
| AI-Q2 | I can skillfully use AI applications or products to help me with my daily work. |
| AI-Q3 | I can choose the most appropriate AI application or product from a variety for a particular task. |
| AI-Q4 | I always comply with ethical principles when using AI applications or products. |
| AI-Q5 | I am never alert to privacy and information security issues when using AI applications or products. |
| AI-Q6 | I am always alert to the abuse of AI technology. |
| AI-Q7 | How frequently do you use generative AI (i.e., artificial intelligence that is capable of producing high quality texts, images, etc. in response to prompts) products such as ChatGPT, Bard, DALL·E 2, Claude, etc.? |
| AI-Q8 | How familiar are you with limitations and shortcomings of generative AI? |

Table 8. AI Literacy Questions. The questions are on a 7 point likert scale ranging from Strongly disagree to Neutral to Strongly agree

| Question ID | Question |
|---|---|
| Util1 | If the only way to save another person's life during an emergency is to sacrifice one's own leg, then one is morally required to make this sacrifice. |
| Util2 | It is morally right to harm an innocent person if harming them is a necessary means to helping several other innocent people. |
| Util3 | From a moral point of view, we should feel obliged to give one of our kidneys to a person with kidney failure since we don't need two kidneys to survive, but really only one to be healthy. |
| Util4 | If the only way to ensure the overall well-being and happiness of the people is through the use of political oppression for a short, limited period, then political oppression should be used. |
| Util5 | From a moral perspective, people should care about the well-being of all human beings on the planet equally; they should not favor the well-being of people who are especially close to them either physically or emotionally |
| Util6 | It is permissible to torture an innocent person if this would be necessary to provide information to prevent a bomb going off that would kill hundreds of people. |
| Util7 | It is just as wrong to fail to help someone as it is to actively harm them yourself. |
| Util8 | Sometimes it is morally necessary for innocent people to die as collateral damage if more people are saved overall. |
| Util9 | It is morally wrong to keep money that one doesn't really need if one can donate it to causes that provide effective help to those who will benefit a great deal. |

Table 9. Utilitarianism Questions. The questions are on a 7-point likert scale ranging from 1 (Strongly Disagree) to 7 (Strongly Agree)

| Question ID | Question |
|---|---|
| Empathy1 | When someone else is feeling excited, I tend to get excited too. |
| Empathy2 | Other people's misfortunes do not disturb me a great deal. |
| Empathy3 | It upsets me to see someone being treated disrespectfully. |
| Empathy4 | I remain unaffected when someone close to me is happy. |
| Empathy5 | I enjoy making other people feel better. |
| Empathy6 | I have tender, concerned feelings for people less fortunate than me. |
| Empathy7 | When a friend starts to talk about his/her problems, I try to steer the conversation towards something else. |
| Empathy8 | I can tell when others are sad even when they do not say anything. |
| Empathy9 | I find that I am "in tune" with other people's moods. |
| Empathy10 | I do not feel sympathy for people who cause their own serious illnesses. |
| Empathy11 | I become irritated when someone cries. |
| Empathy12 | I am not really interested in how other people feel. |
| Empathy13 | I get a strong urge to help when I see someone who is upset. |
| Empathy14 | When I see someone being treated unfairly, I do not feel very much pity for them. |
| Empathy15 | I find it silly for people to cry out of happiness. |
| Empathy16 | When I see someone being taken advantage of, I feel kind of protective towards him/her. |

Table 10. Empathy Questions. The questions are on a 5 point likert scale ranging from Never to Always.

## A.2 Participant Details

The demographics of the participants for our study is shown in Tables 13 to 35. There was a fairly balanced distribution of participants across the different age groups, although there was a slightly higher proportion of participants in the 25-34 years old and 45-54 years old age ranges. In terms of racial distribution, there were more White/Caucasian participants compared to the other races. The gender distribution was relatively balanced in terms of males vs non-males.

| Question ID | Question |
|---|---|
| **Moral Foundation Questionnaire (First Half)** | |
| When you decide whether something is right or wrong, to what extent is the following consideration relevant to your thinking? | |
| MFQ 1 | Whether or not someone suffered emotionally |
| MFQ 2 | Whether or not some people were treated differently than others |
| MFQ 3 | Whether or not someone's action showed love for his or her country |
| MFQ 4 | Whether or not someone showed a lack of respect for authority |
| MFQ 5 | Whether or not someone violated standards of purity and decency |
| MFQ 6 | Whether or not someone was good at math |
| MFQ 7 | Whether or not someone cared for someone weak or vulnerable |
| MFQ 8 | Whether or not someone acted unfairly |
| MFQ 9 | Whether or not someone did something to betray his or her group |
| MFQ 10 | Whether or not someone conformed to the traditions of society |
| MFQ 11 | Whether or not someone did something disgusting |
| **Moral Foundation Questionnaire (Second Half)** | |
| MFQ 12 | Compassion for those who are suffering is the most crucial virtue. |
| MFQ 13 | When the government makes laws, the number one principle should be ensuring that everyone is treated fairly. |
| MFQ 14 | I am proud of my country's history. |
| MFQ 15 | Respect for authority is something all children need to learn. |
| MFQ 16 | People should not do things that are disgusting, even if no one is harmed. |
| MFQ 17 | It is better to do good than to do bad. |
| MFQ 18 | One of the worst things a person could do is hurt a defenseless animal. |
| MFQ 19 | Justice is the most important requirement for a society. |
| MFQ 20 | People should be loyal to their family members, even when they have done something wrong. |
| MFQ 21 | Men and women each have different roles to play in society. |
| MFQ 22 | I would call some acts wrong on the grounds that they are unnatural. |

Table 11. Moral Foundation Questionnaire: 20 Questions. The first part of the questionnaire consists of 11 questions on a 6-point likert scale ranging from 0 (Not At All Relevant) to 5 (Extremely Relevant). The second part of the questionnaire consists of 11 questions on a 6-point likert scale ranging from 0 (Strongly Disagree) to 5 (Strongly Agree). Note: Questions MFQ 6 and 17 are meant to catch participants that are not answering the questionnaire properly and are not included in the MFQ score calculation.

The participants were mostly employed or looking for work and a majority of them also had at least some form of college education. Most participants identified as liberal in terms of political leaning. Participants' AI literacy scores are shown in Table 19 and AI Ethics score are shown in Table 20.

Participants were allocated 5 use cases from one of the scenarios and the allocation between the 2 scenarios are well-balanced and can be found in Table 12.

### A.3 Open-text Annotation Dimensions

*Reasoning Type.* Inspired by previous works in moral psychology, we used two main reasoning types to characterize participants' decision making pattern as expressed in their open-text answers: cost-benefit reasoning and rule-based reasoning [21]. These two reasoning types are rooted in two decision making processes in moral and wider decision making literature: utilitarian and deontological reasoning, respectively. Cost-benefit reasoning thus considers the possible outcomes and their expected utility or value when making decisions, and rule-based reasoning shows more inherent value in action or entities. See § 2.3 for further discussion.

*Moral Values.* To annotate which values were prevalent in participants' considerations of use cases, we used five moral values: care, fairness, loyalty, authority, and purity [33, 34]. While these dimensions have been re-defined to include more diverse values from participants beyond WEIRD (white, educated, industrialized, rich, and democratic) [7], we used these five dimensions due to brevity of the questionnaire, which was used in the survey to provide importance of each values to participants.

| Use Case | Participants Allocated |
|---|---|
| Personal Use Cases | |
| Digital Medical Advice Customized Lifestyle Coach Personal Health Research Nutrition Optimizer Flavorful Swaps | 97 |
| Labor Replacement Use Cases | |
| Lawyer Elementary School Teacher IT Support Specialist Government Eligibility Interviewer Telemarketer | 100 |

Table 12. Participant allocation to each category of scenarios.

| Racial Identity | (N) (%) | Age | N (%) | Gender Identity | N (%) | Education | N (%) |
|---|---|---|---|---|---|---|---|
| White or Caucasian | 33 (33.0) | 18-24 | 11 (11.0) | Man | 49 (49.0) | Bachelor's degree | 36 (36.0) |
| Black or African American | 23 (23.0) | 45-54 | 27 (27.0) | Non-male | 51 (51.0) | Graduate degree* | 18 (18.0) |
| Asian | 21 (21.0) | 25-34 | 22 (22.0) | | | Some college * | 17 (17.0) |
| Mixed | 13 (13.0) | 55-64 | 20 (20.0) | | | High school diploma* | 16 (16.0) |
| Other | 10 (10.0) | 35-44 | 14 (14.0) | | | Associates degree* | 13 (13.0) |
| | 2 (0.7) | 65+ | 6 (6.0) | | | Some high school* | 0 (0.0) |

Table 13. Labor Replacement Study 1 Survey: Racial, age, gender identities and education level of participants. Asterisk (*) denotes labels shortened due to space.

| Minority/Disadvantaged Group | (N) (%) | Transgender | N (%) | Sexuality | N (%) | Political Leaning | N (%) |
|---|---|---|---|---|---|---|---|
| No | 68 (68.0) | No | 97 (97.0) | Heterosexual | 78 (78.0) | Liberal | 34 (34.0) |
| Yes | 32 (32.0) | Yes | 2 (2.0) | Others | 22 (22.0) | Moderate | 23 (23.0) |
| | | Prefer not to say | 1 (1.0) | | | Strongly liberal | 20 (20.0) |
| | | | | | | Conservative | 18 (18.0) |
| | | | | | | Strongly conservative | 4 (4.0) |
| | | | | | | Prefer not to say | 1 (1.0) |

Table 14. Labor Replacement Study 1 Survey: Additional demographic identities

*Switching Conditions.* We annotated concerns expressed in switching conditions using three categories: functionality (e.g., errors, bias in systems, limited capabilities), usage (e.g., errors, bias in systems, limited capabilities), and societal impact (e.g., job loss, over-reliance), inspired by harm taxonomy developed by Solaiman et al. and user concern annotation practice adopted by Mun et al..

## B Open-text Annotation Details

### B.1 Automatic Annotation

*B.1.1 Methods.* We used Open-AI's o1-mini model with maximum tokens set to 1024 to control response length, use a temperature of 0.7 to manage randomness, and keep top_p at 1 with default settings for frequency and presence penalties at 0. Prompts will be released with data upon acceptance.

| Longest Residence | (N) (%) | Employment | N (%) | Occupation (Top 10) | N (%) | Religion | N (%) |
|---|---|---|---|---|---|---|---|
| United States of America | 97 (97.0) | Employed, 40+ | 53 (53.0) | Other | 35 (35.0) | Christian | 29 (29.0) |
| Others | 3 (3.0) | Employed, 1-39 | 16 (16.0) | Prefer not to answer | 10 (10.0) | Agnostic | 20 (20.0) |
| | | Retired | 9 (9.0) | Health Care and Social Assistance | 10 (10.0) | Atheist | 15 (15.0) |
| | | Not employed, looking for work | 7 (7.0) | Information | 10 (10.0) | Nothing in particular | 13 (13.0) |
| | | Disabled, not able to work | 5 (5.0) | Manufacturing | 7 (7.0) | Catholic | 11 (11.0) |
| | | Not employed, NOT looking for work | 4 (4.0) | Professional, Scientific, and Technical Services | 7 (7.0) | Muslim | 5 (5.0) |
| | | Other: please specify | 4 (4.0) | Arts, Entertainment, and Recreation | 6 (6.0) | Hindu | 3 (3.0) |
| | | Prefer not to disclose | 2 (2.0) | Retail Trade | 6 (6.0) | Something else, Specify | 2 (2.0) |
| | | | | Finance and Insurance | 5 (5.0) | Jewish | 1 (1.0) |
| | | | | Transportation and Warehousing, and Utilities | 4 (4.0) | Buddhist | 1 (1.0) |

Table 15. Labor Replacement Study 1 Survey: Additional demographic identities. The Occupation category was capped at the top 10 for brevity, with the remaining occupations merged together with the Other: please specify option.

| Racial Identity | (N) (%) | Age | N (%) | Gender Identity | N (%) | Education | N (%) |
|---|---|---|---|---|---|---|---|
| White or Caucasian | 29 (29.9) | 45-54 | 30 (30.9) | Man | 50 (51.5) | Bachelor's degree | 40 (41.2) |
| Black or African American | 26 (26.8) | 25-34 | 26 (26.5) | Non-male | 47 (48.5) | Some college * | 21 (21.6) |
| Asian | 20 (20.6) | 55-64 | 14 (14.4) | | | Graduate degree* | 14 (14.4) |
| Other | 14 (14.4) | 35-44 | 13 (13.4) | | | High school diploma* | 13 (13.4) |
| Mixed | 8 (8.2) | 18-24 | 9 (9.3) | | | Associates degree* | 8 (8.2) |
| | 2 (0.7) | 65+ | 4 (4.1) | | | Some high school* | 1 (1.0) |
| | | Prefer not to disclose | 1 (1.0) | | | | |

Table 16. Personal Use Cases Study 1 Survey: Racial, age, gender identities and education level of participants. Asterisk (*) denotes labels shortened due to space.

| Minority/Disadvantaged Group | (N) (%) | Transgender | N (%) | Sexuality | N (%) | Political Leaning | N (%) |
|---|---|---|---|---|---|---|---|
| No | 51 (52.6) | No | 94 (96.9) | Heterosexual | 75 (77.3) | Liberal | 34 (35.1) |
| Yes | 46 (47.4) | Yes | 3 (3.1) | Others | 22 (22.7) | Moderate | 31 (32.0) |
| | | Prefer not to say | 0 (0.0) | | | Strongly liberal | 12 (12.4) |
| | | | | | | Conservative | 10 (10.3) |
| | | | | | | Strongly conservative | 9 (9.3) |
| | | | | | | Prefer not to say | 1 (1.0) |

Table 17. Personal Use Cases Study 1 Survey: Additional demographic identities

*B.1.2 Results.* Results for inter-rater reliability analysis of o1's annotations are shown in Table 21.

## C Factors Impacting Acceptability Judgments

### C.1 Use Case Factors

Additional analysis of use case factors showing distribution of judgments by use case sorted by standard deviation is shown in Figure 8. Table 22 shows analysis of use case effect using ANOVA.

| Longest Residence | (N) (%) | Employment | N (%) | Occupation (Top 10) | N (%) | Religion | N (%) |
|---|---|---|---|---|---|---|---|
| United States of America | 93 (95.9) | Employed, 40+ | 46 (47.4) | Other | 36 (35.6) | Christian | 40 (40.8) |
| Others | 4 (4.1) | Employed, 1-39 | 22 (22.7) | Health Care and Social Assistance | 11 (11.3) | Catholic | 16 (16.3) |
| | | Not employed, looking for work | 13 (13.4) | Prefer not to answer | 10 (10.3) | Agnostic | 15 (15.3) |
| | | Not employed, NOT looking for work | 4 (4.1) | Professional, Scientific, and Technical Services | 9 (9.3) | Nothing in particular | 11 (11.2) |
| | | Disabled, not able to work | 4 (4.1) | Educational Services | 9 (9.3) | Atheist | 5 (5.1) |
| | | Other: please specify | 4 (4.1) | Finance and Insurance | 8 (8.2) | Something else, Specify | 5 (5.1) |
| | | Retired | 3 (3.1) | Arts, Entertainment, and Recreation | 5 (5.2) | Buddhist | 3 (3.1) |
| | | Prefer not to disclose | 1 (1.0) | Manufacturing | 5 (5.2) | Muslim | 1 (1.0) |
| | | | | Retail Trade | 4 (4.1) | Jewish | 1 (1.0) |
| | | | | Accommodation and Food Services | 4 (4.1) | Hindu | 1 (1.0) |

Table 18. Personal Use Cases Study 1 Survey: Additional demographic identities. The Occupation category was capped at the top 10 for brevity, with the remaining occupations merged together with the Other: please specify option.

| Score | AI Awareness | AI Usage | AI Evaluation | Gen AI Usage Freq. | Gen AI Limit. Familiarity |
|---|---|---|---|---|---|
| 1 | 25 | 15 | 30 | 35 | 55 |
| 2 | 40 | 60 | 75 | 200 | 320 |
| 3 | 75 | 90 | 105 | 155 | 345 |
| 4 | 140 | 125 | 125 | 235 | 230 |
| 5 | 380 | 310 | 275 | 220 | 35 |
| 6 | 280 | 300 | 310 | 140 | — |
| 7 | 45 | 85 | 65 | — | — |

Table 19. AI literacy scale participant count. Questions are on a 7-point likert scale of increasing score meaning increase in literacy for the aspect. Gen AI Usage Frequency has max score of 6 and Limitation Familiarity has max value of 5.
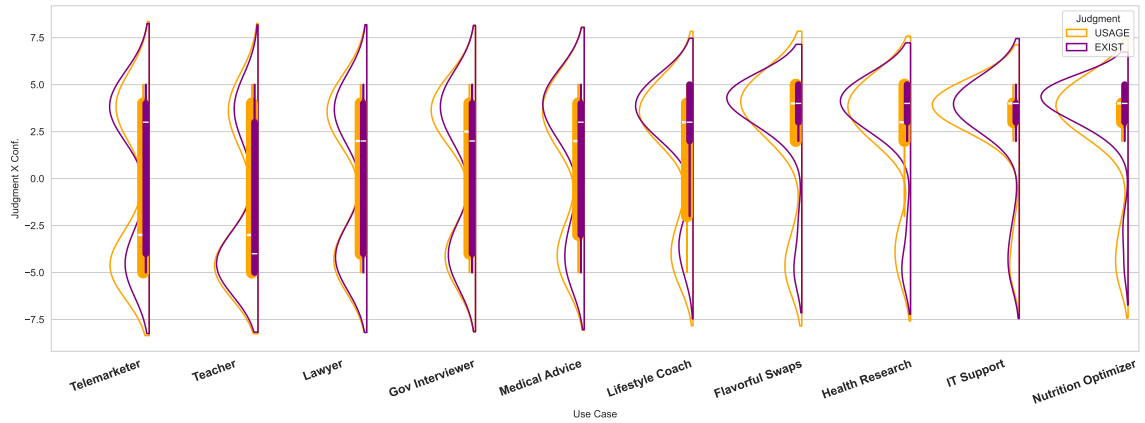


Fig. 8. Numerically converted Judgment x Confidence (-5, 5) by use cases distributions sorted by standard deviation of both existence and usage (sum; highest to lowest) using data from Study 1 results.

| Score | AI Ethics |
|---|---|
| 5 | 15 |
| 6 | 10 |
| 7 | 25 |
| 8 | 25 |
| 9 | 75 |
| 10 | 105 |
| 11 | 165 |
| 12 | 100 |
| 13 | 105 |
| 14 | 90 |
| 15 | 120 |
| 16 | 80 |
| 17 | 35 |
| 18 | 35 |

Table 20. AI ethics score count for total AI ethics score (sum over 3 questions with 7 point likert scale with max possible value of 21)

| Category | AC1 | Interpretation | 95% CI | p-value | z | SE | PA | PE |
|---|---|---|---|---|---|---|---|---|
| Cost Benefit | 0.976 | Almost Perfect | (0.942, 1.000) | **0.0** | 56.9 | 0.0172 | 0.980 | 0.164 |
| Rule Based | 0.848 | Almost Perfect | (0.754, 0.943) | **0.0** | 17.8 | 0.0476 | 0.890 | 0.276 |
| Care | 0.427 | Moderate | (0.234, 0.619) | $2.76 \times 10^{-5}$ | 4.40 | 0.0970 | 0.650 | 0.390 |
| Fairness | 0.411 | Moderate | (0.224, 0.599) | $3.29 \times 10^{-5}$ | 4.35 | 0.0945 | 0.690 | 0.474 |
| Authority | 0.839 | Almost Perfect | (0.743, 0.935) | **0.0** | 17.3 | 0.0486 | 0.880 | 0.255 |
| Purity | 0.758 | Substantial | (0.639, 0.878) | **0.0** | 12.6 | 0.0603 | 0.820 | 0.255 |
| Functionality | 0.587 | Moderate | (0.425, 0.749) | $1.31 \times 10^{-10}$ | 7.18 | 0.0818 | 0.790 | 0.492 |
| Usage | 0.673 | Substantial | (0.525, 0.822) | $1.47 \times 10^{-14}$ | 9.03 | 0.0746 | 0.820 | 0.449 |
| Societal Impact | 0.595 | Moderate | (0.432, 0.759) | $1.01 \times 10^{-10}$ | 7.23 | 0.0823 | 0.770 | 0.432 |

Table 21. Inter-rater Agreement using Gwet's AC1. Interpretation according to [81].

| Acceptability | Aspect | Factor | Sum Sq | Mean Sq | NumDF | DenDF | Pr(>F) |
|---|---|---|---|---|---|---|---|
| EXIST | Judgment | Category | 29.903 | 29.903 | 1 | 197 | **5.98e-11** *** |
| | | Use Case | 86.116 | 9.5684 | 9 | 641.38 | **< 2.2e-16** *** |
| | Confidence | Category | 0.0017113 | 0.0017113 | 1 | 197 | 0.9563 |
| | | Use Case | 13.519 | 1.5021 | 9 | 603.23 2.7243 | **0.004037** ** |
| USAGE | Judgment | Category | 8.4257 | 8.4257 | 1 | 197 | **0.0002488** *** |
| | | Use Case | 73.801 | 8.2001 | 9 | 610.34 | **< 2.2e-16** *** |
| | Confidence | Category | 1.3444 | 1.3444 | 1 | 197 | 0.153 |
| | | Use Case | 20.721 | 2.3023 | 9 | 603.36 | **0.0001783** *** |

Table 22. ANOVA analysis of LMER models judgment ~ category + (1 | subject) and judgment ~ useCase + (1 | subject) (same formulas repeated with confidence as a dependent variable) analyzed with Study 1 data.

## C.2 Demographics Factors

Additional analysis using ANOVA for demographic factors is shown in Table 23.

*C.2.1 Questionnaires.* Interestingly, only Loyalty had a significant effect on both existence (0.20, $p < .001$) and usage (0.20, $p < .01$) as shown in Table 24. Moreover, Empathy had a positive and marginally significant effect for usage

| Demographics | EXIST | | | USAGE | | |
| --- | --- | --- | --- | --- | --- | --- |
| | Judgment | Confidence | Judg.×Conf. | Judgment | Confidence | Judg×Conf. |
| **All** | | | | | | |
|   Gender | **16.60**\*** | 0.19 | **16.71**\*** | **15.26**\*** | 0.83 | **13.14**\*** |
|   Race | 1.62 | **4.09**\** | 1.45 | 0.65 | **5.12**\*** | 0.45 |
|   Employment | 1.13 | **3.03**\* | 1.14 | 0.43 | 1.71 | 0.69 |
|   Sexual Orientation | 0.42 | 0.37 | 0.19 | 0.75 | **5.22**\* | 0.09 |
| **Professional** | | | | | | |
|   Race | **2.56**\* | 1.80 | 2.34· | 1.04 | **2.91**\* | 0.48 |
|   Gender | **18.37**\*** | 0.05 | **19.51**\*** | **19.83**\*** | 0.19 | **20.21**\*** |
|   Education | 1.98 | 0.96 | 1.34 | 2.25· | 1.12 | 2.07· |
|   Discrimination | 2.18 | 0.68 | 2.67· | 0.29 | 1.46 | 0.13 |
| **Personal** | | | | | | |
|   Race | 2.11· | **4.36**\** | 2.28· | 0.16 | **4.07**\** | 0.38 |
|   Political View | 0.38 | **3.39**\* | 0.86 | 0.33 | 1.56 | 0.36 |
|   Employment | 0.85 | **2.42**\* | 1.47 | 0.33 | 0.30 | 0.36 |

\*\*\*$p < 0.001$; \*\*$p < 0.01$; \*$p < 0.05$; ·$p < 0.1$

Table 23. ANOVA Results by Demographic Category (F-value with Significance)

| Decision Style Factors | EXIST ($\beta$ (SE)) | | | USAGE ($\beta$ (SE)) | | |
| --- | --- | --- | --- | --- | --- | --- |
| | Judg. | Conf. | Judg.×Conf. | Judg. | Conf. | Judg.×Conf. |
| (Intercept) | 0.11 (0.34) | **3.10**\*** (0.46) | −0.32 (1.46) | −0.74 (0.39) | **2.98**\*** (0.49) | **−4.04**\* (1.68) |
| MFQ Care | 0.00 (0.01) | −0.01 (0.02) | 0.01 (0.06) | −0.00 (0.02) | 0.02 (0.02) | −0.01 (0.07) |
| MFQ Fairness | −0.01 (0.02) | 0.04 (0.02) | −0.03 (0.07) | 0.01 (0.02) | 0.01 (0.02) | 0.05 (0.08) |
| MFQ Loyalty | **0.04**\*** (0.01) | −0.01 (0.02) | **0.18**\*** (0.05) | **0.04**\** (0.01) | −0.01 (0.02) | **0.18**\** (0.06) |
| MFQ Authority | −0.01 (0.01) | 0.02 (0.02) | −0.05 (0.05) | 0.00 (0.02) | 0.03 (0.02) | 0.02 (0.06) |
| MFQ Purity | 0.00 (0.01) | 0.02 (0.01) | 0.04 (0.04) | −0.01 (0.01) | 0.00 (0.02) | −0.01 (0.05) |
| Empathy | 0.01 (0.03) | 0.03 (0.04) | 0.05 (0.13) | 0.06 (0.04) | 0.02 (0.05) | 0.28 (0.15) |
| InstrumentalHarm | 0.00 (0.03) | −0.01 (0.04) | 0.04 (0.13) | −0.02 (0.03) | −0.01 (0.04) | −0.06 (0.15) |
| ImpartialBenificence | 0.03 (0.03) | −0.01 (0.04) | 0.08 (0.13) | 0.02 (0.04) | −0.02 (0.05) | 0.06 (0.15) |
| AIC | 2412.17 | 2539.14 | 5149.55 | 2442.29 | 2671.82 | 5132.26 |
| BIC | 2470.88 | 2597.85 | 5208.26 | 2501.01 | 2730.53 | 5190.97 |
| Log Likelihood | −1194.08 | −1257.57 | −2562.77 | −1209.15 | −1323.91 | −2554.13 |
| Num. obs. | 985 | 985 | 985 | 985 | 985 | 985 |
| Num. groups: prolific_id | 197 | 197 | 197 | 197 | 197 | 197 |
| Num. groups: use_case | 10 | 10 | 10 | 10 | 10 | 10 |
| Var: prolific_id (Intercept) | 0.12 | 0.37 | 2.71 | 0.23 | 0.42 | 4.65 |
| Var: use_case (Intercept) | 0.11 | 0.01 | 2.22 | 0.08 | 0.02 | 1.58 |
| Var: Residual | 0.55 | 0.56 | 8.53 | 0.53 | 0.64 | 7.70 |

\*\*\*$p < 0.001$; \*\*$p < 0.01$; \*$p < 0.05$

Table 24. Coefficients with standard error in parenthesis with following models: Judgment ~ MFQ$_{foundation}$ + Empathy + InustrumentalHarm + ImpartialBeneficence + (1|Subject) + (1|useCase). Bolded value for empathy had $p < 0.1$.

$(.09, p < .1)$. However, Loyalty, as shown in Figure 6, does not appear as frequently in participants' open text responses compared to values such as Care and Fairness and is the only value that did not have a significant association with use cases.

## D  Factors in Participant Rationale

### D.1  Reasoning Types

We show the flow of participants' decisions and corresponding rationales throughout use cases in Figure 9, which shows interesting distribution and switching of reasoning types, which would be interesting for future studies to consider. Moreover, Table 25 shows that there are almost no relation between reasoning types used by the participants and the decision-making style questionnaire results signifying that the reasoning types might be highly use-case specific rather than a character trait. It would be interesting to study the factors that actually influence the choice of reasoning types.



(a)  Acceptance judgment and reasoning type mapping throughout professional use cases.



(b)  Acceptance judgment and reasoning type mapping throughout personal use cases.

Fig. 9.  Mapping of decisions and reasoning types. + and - denote positive and negative acceptance. "C" denotes Cost-benefit analysis and "R" denotes Rule-based reasoning.

| | Cost-benefit | Rule-based |
|---|---|---|
| (Intercept) | **0.73**\*\*\*(0.13) | **0.28**\*(0.13) |
| MFQ Care | −0.00(0.01) | 0.00(0.01) |
| MFQ Fairness | 0.01(0.01) | 0.00(0.01) |
| MFQ Loyalty | **0.01**\*(0.00) | −0.01(0.01) |
| MFQ Authority | −0.00(0.01) | 0.00(0.01) |
| MFQ Purity | −0.00(0.00) | 0.00(0.00) |
| Empathy | 0.00(0.01) | −0.01(0.01) |
| InstrumentalHarm | 0.00(0.01) | −0.00(0.01) |
| ImpartialBenificence | −0.00(0.01) | −0.00(0.01) |
| AIC | 668.43 | 815.32 |
| BIC | 727.14 | 874.03 |
| Log Likelihood | −322.21 | −395.66 |
| Num. obs. | 985 | 985 |
| Num. groups: prolific_id | 197 | 197 |
| Num. groups: use_case | 10 | 10 |
| Var: prolific_id (Intercept) | 0.02 | 0.02 |
| Var: use_case (Intercept) | 0.00 | 0.01 |
| Var: Residual | 0.10 | 0.11 |

$^{***}p < 0.001; ^{**}p < 0.01; ^{*}p < 0.05$

Table 25. Coefficient and standard error with significance. Model defined by reasoningType ~ MFQ$_{foundation}$ + Empathy + InstrumentalHarm + ImpartialBeneficence + (1|subject) + (1|useCase)

## D.2    Impact of Rationale Factors on Judgment

We display the analysis results using ANOVA to understand the effect of rationale factors on judgment in Table 26.

## D.3    Factors Influencing Moral Foundations in Rationale

In Table 27 we display analysis result using linear mixed effects on factors that may influence moral foundations appealed to in participants' rationales. We find greater relations with the use cases than personal factors.

## E    Survey 2: Explicit Weighing of Harms and Benefits of Use Cases

Although not included in main text, we administered a variation of our main study where we asked participants to explicitly reason through harms and benefits. The decisions were measured before and after the explicit weighing of harms and benefits. However, we saw almost no effect.

### E.1    Study Overview

To better understand the reasoning behind participants decisions about the judgment and usage of the use cases, we conducted a second study with 1 survey for each category (Labor Replacement Use Cases and Personal Use Cases). The second study includes an additional set of questions to elicit the harms and benefits of developing and not developing an use case to better understand the reasoning behind participants decisions. Furthermore, we asked participants the judgment and usage decision questions before and after the set of harms and benefits questions to see if listing out reasons about an use case would elicit any change in their decisions. The details of the second study can be found in E.2 and the results can be found in E.3.

| Acceptability | Factor | Sum Sq | Mean Sq | NumDF | DenDF | Pr(>F) |
|---|---|---|---|---|---|---|
| | **Judgment** | | | | | |
| | Cost Benefit | 2.941 | 2.941 | 1 | 955.32 | **0.001306** ** |
| | Rule Based | 7.187 | 7.187 | 1 | 941.74 | **5.562e-07** *** |
| | Fairness | 1.551 | 1.551 | 1 | 966.41 | **0.019414** * |
| | Authority | 1.044 | 1.044 | 1 | 961.50 | 0.055021 . |
| | Functionality | 1.137 | 1.137 | 1 | 964.09 | **0.045230** * |
| | Usage | 125.594 | 125.594 | 1 | 837.10 | **< 2.2e-16** *** |
| **EXIST** | Societal Impact | 1.873 | 1.873 | 1 | 967.89 | **0.010228** * |
| | **Confidence** | | | | | |
| | Rule Based | 5.7805 | 5.7805 | 1 | 856.38 | **0.001158** ** |
| | Care | 3.8723 | 3.8723 | 1 | 907.71 | **0.007762** ** |
| | Fairness | 5.7117 | 5.7117 | 1 | 899.21 | **0.001237** ** |
| | Authority | 2.2665 | 2.2665 | 1 | 851.91 | **0.041526** * |
| | Usage | 2.6254 | 2.6254 | 1 | 809.17 | **0.028303** * |
| | **Judgment x Confidence** | | | | | |
| | Cost Benefit | 48.04 | 48.04 | 1 | 940.91 | **0.0004985** *** |
| | Rule Based | 86.12 | 86.12 | 1 | 924.15 | **3.331e-06** *** |
| | Fairness | 31.87 | 31.87 | 1 | 972.26 | **0.0045236** ** |
| | Authority | 14.24 | 14.24 | 1 | 972.97 | 0.0574858 . |
| | Functionality | 18.99 | 18.99 | 1 | 971.66 | **0.0283005** * |
| | Usage | 2454.30 | 2454.30 | 1 | 888.43 | **< 2.2e-16** *** |
| | Societal Impact | 34.26 | 34.26 | 1 | 971.98 | **0.0032486** ** |
| | **Judgment** | | | | | |
| | Cost Benefit | 0.28 | 0.28 | 1 | 933.91 | **0.01742** * |
| | Usage | 491.35 | 491.35 | 1 | 943.16 | **< 2e-16** *** |
| | **Confidence** | | | | | |
| | Cost Benefit | 2.2380 | 2.2380 | 1 | 864.41 | 0.0555024 . |
| | Rule Based | 2.0230 | 2.0230 | 1 | 852.41 | 0.0686374 . |
| **USAGE** | Fairness | 4.2517 | 4.2517 | 1 | 895.25 | **0.0083622** ** |
| | Authority | 2.4525 | 2.4525 | 1 | 852.20 | **0.0450309** * |
| | Loyalty | 2.2751 | 2.2751 | 1 | 875.10 | 0.0535163 . |
| | Functionality | 2.6366 | 2.6366 | 1 | 897.34 | **0.0376892** * |
| | Usage | 9.2673 | 9.2673 | 1 | 817.52 | **0.0001033** *** |
| | Societal Impact | 2.1237 | 2.1237 | 1 | 913.41 | 0.0620930 . |
| | **Judgment x Confidence** | | | | | |
| | Cost Benefit | 106.358 | 106.358 | 1 | 874.26 | **8.333e-05** *** |
| | Rule Based | 33.512 | 33.512 | 1 | 859.58 | **0.026742** * |
| | Fairness | 89.120 | 89.120 | 1 | 930.00 | **0.000312** *** |
| | Societal Impact | 74.975 | 74.975 | 1 | 929.00 | **0.000938** *** |

Table 26. ANOVA analysis of the LMER model results.

## E.2  Setup and Details

While these same set of questions are asked for all five use cases for our main study, in our second study, participants are randomly allocated a single use case. The second study differs from the main study with an initial set of judgment questions without open-text rationales (Q1 - Initial to Q4 - Initial), which are followed by explicit listing and weighing of the possible harms and benefits of the use case in the context of both developing and not developing the use case. We then again ask participants the same set of judgment questions along with the open-text questions to elaborate on their reasoning, similar to the main study. To understand how the judgment and usage decisions are affected by other factors,

| | Care | Fairness | Purity | Loyalty | Authority |
|---|---|---|---|---|---|
| (Intercept) (Telemarketer) | **1.74***** **(0.31)** | 0.29 (0.23) | **−2.80***** **(0.42)** | **−4.82***** **(1.03)** | **−2.29***** **(0.35)** |
| MFQ_care | 0.08 (0.19) | −0.12 (0.13) | −0.36 (0.19) | −0.15 (0.43) | −0.01 (0.19) |
| MFQ_fairness | 0.13 (0.19) | **0.31*** **(0.13)** | 0.30 (0.20) | −0.19 (0.43) | 0.14 (0.19) |
| MFQ_loyalty | **0.60**** **(0.21)** | **0.31*** **(0.14)** | −0.14 (0.21) | −0.34 (0.57) | −0.12 (0.20) |
| MFQ_authority | **−0.48*** **(0.24)** | −0.23 (0.16) | 0.23 (0.24) | 0.06 (0.56) | 0.23 (0.24) |
| MFQ_purity | 0.28 (0.19) | −0.16 (0.13) | −0.07 (0.20) | −0.05 (0.44) | −0.22 (0.19) |
| empathy_total | 0.06 (0.15) | 0.04 (0.10) | 0.07 (0.15) | 0.24 (0.39) | −0.05 (0.15) |
| InstrumentalHarm | −0.04 (0.15) | 0.12 (0.10) | −0.14 (0.16) | −0.10 (0.39) | 0.09 (0.15) |
| ImpartialBenificence | 0.09 (0.15) | −0.05 (0.10) | **−0.37*** **(0.16)** | −0.44 (0.41) | −0.07 (0.15) |
| Gov. Eligi. Interviewer | 0.06 (0.39) | **1.45***** **(0.35)** | **−1.72*** **(0.82)** | — | −0.01 (0.44) |
| IT Support Specialist | **1.14*** **(0.46)** | **0.76*** **(0.32)** | 0.44 (0.49) | — | −0.72 (0.50) |
| Elementary School Teacher | **0.88*** **(0.44)** | **−0.66*** **(0.31)** | 0.90 (0.47) | 1.43 (1.13) | 0.42 (0.42) |
| Lawyer | −0.48 (0.37) | **0.70*** **(0.32)** | −0.71 (0.60) | 0.71 (1.24) | **1.41***** **(0.40)** |
| Flavorful Swaps | 0.73 (0.47) | −0.18 (0.33) | **1.50**** **(0.48)** | −0.04 (1.43) | −0.80 (0.55) |
| Nutrition Optimizer | **1.14*** **(0.50)** | −0.37 (0.33) | −0.01 (0.55) | — | −0.26 (0.50) |
| Personal Health Research | **1.46**** **(0.54)** | 0.53 (0.34) | −0.32 (0.58) | — | 0.42 (0.47) |
| Cust, Lifestyle Coach | 0.71 (0.47) | 0.47 (0.34) | 0.56 (0.51) | −0.04 (1.43) | −0.64 (0.54) |
| Digital Medical Advice | **1.45**** **(0.53)** | 0.02 (0.33) | −0.50 (0.60) | — | **1.47**** **(0.44)** |
| AIC | 750.75 | 1256.07 | 644.21 | 120.53 | 845.93 |
| BIC | 843.71 | 1349.03 | 737.17 | 213.49 | 938.89 |
| Log Likelihood | −356.38 | −609.04 | −303.11 | −41.26 | −403.97 |
| Num. obs. | 985 | 985 | 985 | 985 | 985 |
| Num. groups: prolific_id | 197 | 197 | 197 | 197 | 197 |
| Var: prolific_id (Intercept) | 1.27 | 0.63 | 1.10 | 0.00 | 1.56 |

$^{***}p < 0.001; {}^{**}p < 0.01; {}^{*}p < 0.05$

Table 27. Effects and standard error in parenthesis of the annotation output of participant answers modeled with following formula $\text{Annot}_{foundation} \sim \text{MFQ}_{foundation}$ + empathy + instrumentalHarm + impartialBeneficence + useCase + (1|subject) using glmer with family set to binomial. Intercept shows effects when categorical variables are set to following: useCase = Telemarketer and Type = Cost-Benefit.

we asked the participants about their demographics, ai literacy levels and several other reasoning factors after the main set of questions, and these questions can be found in §A.1 The main questions for the second study can be found in Table 29. The participant demographics for the second study can be found in Tables 30 to 35. The distribution of each use case within each scenario (Labor Replacement Use Cases and Personal Use Cases) for the second study is relatively well-balanced and can be found in Table 28.

### E.3 Results

To explore the possible impact of explicitly weighing harms and benefits of a use case on participant's decision, we analyzed the participant's judgment of acceptability before and after explicit weighing of harms and benefits (Study 2; see § ?? for details on questions asked). The Type III ANOVA with Satterthwaite's method for measurement time (before, after) indicated a marginally significant effect $F(1, 201.05) = 3.371, p = 0.0678$ on usage judgment weighed by confidence, which suggests that explicit harms and benefits weighing may have an influence, albeit not at conventional significance levels. We further analyzed reasoning effect on each subset of data pertinent to each use case through a mixed effects regression model with judgment metric as a dependent variable and measurement time as an independent variable with random effect from subject. Interestingly, the result was significant for Customized Lifestyle Coach AI across different judgments including, existence ($\beta = -0.40, SE = 0.18, p < .05$), confidence-weighed

| Use Case | Participants Allocated |
|---|---|
| Personal Use Cases | |
| Digital Medical Advice | 20 |
| Customized Lifestyle Coach | 20 |
| Personal Health Research | 19 |
| Nutrition Optimizer | 21 |
| Flavorful Swaps | 17 |
| Labor Replacement Use Cases | |
| Lawyer | 20 |
| Elementary School Teacher | 22 |
| IT Support Specialist | 17 |
| Government Eligibility Interviewer | 19 |
| Telemarketer | 23 |

Table 28. Use Case allocation for Study 2. Specific participant numbers are listed for each use case.

existence ($\beta = -1.05, SE = 0.51, p < .05$), and confidence-weighed usage judgments ($\beta = -0.75, SE = 0.38, p < .05$). Explicit weighing also had a significant effect on confidence of existence judgment for Digital Medical Advice AI ($\beta = 0.30, SE = 0.11, p < .01$). The negative coefficients for Customized Lifestyle AI suggests that weighing harms and benefits caused participants to lower acceptance and positive coefficient to confidence on judgments on Digital Medical AI suggests that weighing harms and benefits solidified decisions. These diverging effects signify an interesting interaction between use cases and explicit weighing of harms and benefits.

| Question ID | Question | Answer Type |
|---|---|---|
| **AI Perception Question (Before)** | | |
| AI Perception Before | Overall, how does the growing presence of artificial intelligence (AI) in daily life and society make you feel? | 5 Point Likert Scale |
| **Initial Decision/Usage** | | |
| Q1 - Initial | Do you think a technology like this should be developed? | Yes/No |
| Q2 - Initial | How confident are you in your above answer? | 5 Point Likert Scale |
| Q3 - Initial | If [Use Case] exists, would you ever use its services (answer yes, even if you think you would use it very infrequently)? | Yes/No |
| Q4 - Initial | How confident are you in your above answer? | 5 Point Likert Scale |
| **Benefits of Developing Use Case** | | |
| Q1 - BDev | How will [Use Case] positively impact individuals? | Text |
| Q2 - BDev | Which groups of people do you think would benefit the most from the above positive impacts? (You can list more than one group.) | Text |
| Q3 - BDev | How beneficial would [Use Case] be if it had the above positive impacts? | 9 Point Likert Scale |
| **Malicious Uses of Developing Use Case** | | |
| Q1 - HDev | Please complete the following: [Use Case] could have a negative impact if it was used to... | Text |
| Q1 - HDev | What would be the negative impact of the above malicious or unintended uses? | Text |
| Q2 - HDev | Which groups of people do you think would be harmed the most by the above malicious or unintended uses? (You can list more than one group.) | Text |
| Q3 - HDev | How harmful would [Use Case] be if it had the above negative impacts? | 9 Point Likert Scale |
| **Failures of Developing Use Case** | | |
| Q1 - HDevF | Please complete the following: If [Use Case] failed to do its intended task properly, fully, and accurately, it could have a negative impact if it... | Text |
| Q1 - HDevF | What would be the negative impact of those failure cases? | Text |
| Q2 - HDevF | Which groups of people do you think would be harmed the most by the above failure cases? (You can list more than one group.) | Text |
| Q3 - HDevF | How harmful would [Use Case] be if it had the above negative impacts? | 9 Point Likert Scale |
| **Benefits of Not Developing Use Case** | | |
| Q1 - BNonDev | Please complete the following: Not having [Use Case] would be beneficial because... | Text |
| Q2 - BNonDev | Which groups of people do you think would benefit the most by banning or not developing [Use Case]? (You can list more than one group.) | Text |
| Q3 - BNonDev | How beneficial would it be if [Use Case] was banned or not developed and it had the above positive impact? | 9 Point Likert Scale |
| **Harms of Not Developing Use Case** | | |
| Q1 - HNonDev | Please complete the following: Not having [Use Case] would be harmful because... | Text |
| Q2 - HNonDev | Which groups of people do you think would be harmed the most by banning or not developing [Use Case]? (You can list more than one group.) | Text |
| Q3 - HNonDev | How harmful would it be if [Use Case] was banned or not developed and it had the above negative impacts? | 9 Point Likert Scale |
| **Final Decision/Usage** | | |
| Q1 - Final | Do you think a technology like this should be developed? | Yes/No |
| Q2 - Final | How confident are you in your above answer? | 5 Point Likert Scale |
| Q3 - Final - Y | Please elaborate on your answer to the previous question: Do you think a technology like this should be developed?: [Q1 - Final Answer] | Text |
| Q3 - Final - N | Please elaborate on your answer to the previous question: Do you think a technology like this should be developed?: [Q1 - Final Answer] | Text |
| Q4 - Final - Y | Under what circumstances would you switch your decision from [Q1 - Final Answer] should be developed to should not be developed? | Text |
| Q4 - Final - N | Under what circumstances would you switch your decision from [Q1 - Final Answer] should not be developed to should be developed? | Text |
| Q5 - Final | If [Use Case] exists, would you ever use its services (answer yes, even if you think you would use it very infrequently)? | Yes/No |
| Q6 - Final | How confident are you in your above answer? | 5 Point Likert Scale |
| **AI Perception Question (After)** | | |
| AI Perception After | Before we continue, we'd like to get your thoughts on AI one more time. Overall, how does the growing presence of artificial intelligence (AI) in daily life and society make you feel? | 5 Point Likert Scale |

Table 29. Study 2 Specific Question. The placeholder [Use Case] is used in place of the 10 use cases chosen for the studies.

| Racial Identity | (N) (%) | Age | N (%) | Gender Identity | N (%) | Education | N (%) |
|---|---|---|---|---|---|---|---|
| White or Caucasian | 32 (31.4) | 45-54 | 32 (31.4) | Man | 49 (49.0) | Bachelor's degree | 44 (43.1) |
| Black or African American | 25 (24.5) | 25-34 | 29 (28.4) | Non-male | 51 (51.0) | Graduate degree* | 18 (17.6) |
| Asian | 18 (17.6) | 35-44 | 12 (11.8) | | | Some college * | 18 (17.6) |
| Mixed | 15 (14.7) | 55-64 | 12 (11.8) | | | High school diploma* | 15 (14.7) |
| Other | 12 (11.8) | 18-24 | 9 (8.8) | | | Associates degree* | 7 (6.9) |
| | 2 (0.7) | 65+ | 8 (7.8) | | | Some high school* | 0 (0.0) |

Table 30. Labor Replacement Study 2 Survey: Racial, age, gender identities and education level of participants. Asterisk (*) denotes labels shortened due to space.

| Minority/Disadvantaged Group | (N) (%) | Transgender | N (%) | Sexuality | N (%) | Political Leaning | N (%) |
|---|---|---|---|---|---|---|---|
| No | 56 (54.9) | No | 97 (95.1) | Heterosexual | 76 (74.5) | Liberal | 37 (36.3) |
| Yes | 46 (45.1) | Yes | 4 (3.9) | Others | 26 (25.5) | Moderate | 27 (26.5) |
| | | Prefer not to say | 1 (1.0) | | | Conservative | 17 (16.7) |
| | | | | | | Strongly liberal | 16 (15.7) |
| | | | | | | Strongly conservative | 4 (3.9) |
| | | | | | | Prefer not to say | 1 (1.0) |

Table 31. Labor Replacement Study 2 Survey: Additional demographic identities

| Longest Residence | (N) (%) | Employment | N (%) | Occupation (Top 10) | N (%) | Religion | N (%) |
|---|---|---|---|---|---|---|---|
| United States of America | 96 (94.1) | Employed, 40+ | 44 (43.1) | Other | 34 (33.3) | Christian | 38 (37.3) |
| Others | 6 (5.9) | Employed, 1-39 | 28 (27.5) | Educational Services | 11 (10.8) | Agnostic | 19 (18.6) |
| | | Not employed, looking for work | 19 (18.6) | Health Care and Social Assistance | 10 (10.0) | Catholic | 18 (17.6) |
| | | Retired | 4 (3.9) | Information | 8 (7.8) | Nothing in particular | 12 (11.8) |
| | | Not employed, NOT looking for work | 3 (2.9) | Prefer not to answer | 8 (7.8) | Atheist | 7 (6.9) |
| | | Other: please specify | 3 (2.9) | Retail Trade | 7 (6.9) | Muslim | 3 (2.9) |
| | | Disabled, not able to work | 1 (1.0) | Finance and Insurance | 7 (6.9) | Something else, Specify | 3 (2.9) |
| | | Prefer not to disclose | 0 (0.0) | Professional, Scientific, and Technical Services | 6 (5.9) | Jewish | 1 (1.0) |
| | | | | Manufacturing | 6 (5.9) | Hindu | 1 (1.0) |
| | | | | Administrative and support and waste management services | 5 (4.9) | Buddhist | 0 (0.0) |

Table 32. Labor Replacement Study 2 Survey: Additional demographic identities. The Occupation category was capped at the top 10 for brevity, with the remaining occupations merged together with the Other: please specify option.

| Racial Identity | (N) (%) | Age | N (%) | Gender Identity | N (%) | Education | N (%) |
|---|---|---|---|---|---|---|---|
| White or Caucasian | 35 (36.1) | 45-54 | 29 (29.9) | Non-male | 50 (51.5) | Bachelor's degree | 33 (34.0) |
| Black or African American | 22 (22.7) | 25-34 | 22 (22.7) | Man | 47 (48.5) | Graduate degree* | 24 (24.7) |
| Asian | 19 (19.6) | 55-64 | 17 (17.5) | | | Some college * | 18 (18.6) |
| Mixed | 13 (13.4) | 35-44 | 17 (17.5) | | | High school diploma* | 12 (12.4) |
| Other | 8 (8.2) | 18-24 | 6 (6.2) | | | Associates degree* | 10 (10.3) |
| | 2 (0.7) | 65+ | 6 (6.2) | | | Some high school* | 0 (0.0) |
| | | Prefer not to disclose | 0 (0.0) | | | | |

Table 33. Personal Use Cases Study 2 Survey: Racial, age, gender identities and education level of participants. Asterisk (*) denotes labels shortened due to space.

| Minority/Disadvantaged Group | (N) (%) | Transgender | N (%) | Sexuality | N (%) | Political Leaning | N (%) |
|---|---|---|---|---|---|---|---|
| No | 50 (51.5) | No | 93 (95.9) | Heterosexual | 76 (78.4) | Liberal | 34 (35.1) |
| Yes | 47 (48.5) | Yes | 4 (4.1) | Others | 21 (21.6) | Moderate | 26 (26.8) |
| | | Prefer not to say | 0 (0.0) | | | Strongly liberal | 17 (17.5) |
| | | | | | | Conservative | 13 (13.4) |
| | | | | | | Strongly conservative | 6 (6.2) |
| | | | | | | Prefer not to say | 1 (1.0) |

Table 34. Personal Use Cases Study 2 Survey: Additional demographic identities

| Longest Residence | (N) (%) | Employment | N (%) | Occupation (Top 10) | N (%) | Religion | N (%) |
|---|---|---|---|---|---|---|---|
| United States of America | 95 (97.9) | Employed, 40+ | 44 (45.4) | Other | 36 (37.1) | Christian | 43 (44.3) |
| Others | 2 (2.1) | Employed, 1-39 | 23 (23.7) | Health Care and Social Assistance | 13 (13.4) | Agnostic | 12 (12.4) |
| | | Not employed, looking for work | 9 (9.3) | Information | 9 (9.3) | Atheist | 12 (12.4) |
| | | Other: please specify | 7 (7.2) | Finance and Insurance | 8 (8.2) | Catholic | 10 (10.3) |
| | | Retired | 6 (6.2) | Prefer not to answer | 6 (6.2) | Nothing in particular | 8 (8.2) |
| | | Disabled, not able to work | 5 (5.2) | Retail Trade | 6 (6.2) | Muslim | 4 (4.1) |
| | | Not employed, NOT looking for work | 2 (2.1) | Manufacturing | 5 (5.1) | Something else, Specify | 4 (4.1) |
| | | Prefer not to disclose | 1 (1.0) | Educational Services | 5 (5.1) | Buddhist | 2 (2.1) |
| | | | | Arts, Entertainment, and Recreation | 5 (5.1) | Hindu | 1 (1.0) |
| | | | | Accommodation and Food Services | 4 (4.1) | Jewish | 1 (1.0) |

Table 35. Personal Use Cases Study 1 Survey: Additional demographic identities. The Occupation category was capped at the top 10 for brevity, with the remaining occupations merged together with the Other: please specify option.