# Organize, Then Vote: Exploring Cognitive Load in Quadratic Survey Interfaces

Ti-Chung Cheng
Computer Science
University of Illinois at
Urbana-Champaign
Urbana, Illinois, USA
tcheng10@illinois.edu

Yutong Zhang*
Computer Science
Stanford University
Stanford, California, USA
yutongz7@stanford.edu

Yi-Hung Chou*
University of California,
Irvine
Irvine, California, USA
yihungc1@uci.edu

Vinay Koshy
Computer Science
University of Illinois at
Urbana Champaign
Urbana, Illinois, USA
vkoshy2@illinois.edu

Tiffany Wenting Li
Computer Science
University of Illinois at
Urbana-Champaign
Urbana, Illinois, USA
wenting7@illinois.edu

Karrie Karahalios
Computer Science
University of Illinois at
Urbana-Champaign
Urbana, Illinois, USA
kkarahal@illinois.edu

Hari Sundaram
Computer Science
University of Illinois
Urbana, Illinois, USA
hs1@illinois.edu

## Abstract

Quadratic Surveys (QSs) elicit more accurate preferences than traditional methods like Likert-scale surveys. However, the cognitive load associated with QSs has hindered their adoption in digital surveys for collective decision-making. We introduce a two-phase "organize-then-vote" QS to reduce cognitive load. As interface design significantly impacts survey results and accuracy, our design scaffolds survey takers' decision-making while managing the cognitive load imposed by QS. In a 2x2 between-subject in-lab study on public resource allotment, we compared our interface with a traditional text interface across a QS with 6 (short) and 24 (long) options. Two-phase interface participants spent more time per option and exhibited shorter voting edit distances. We qualitatively observed shifts in cognitive effort from mechanical operations to constructing more comprehensive preferences. We conclude that this interface promoted deeper engagement, potentially reducing satisficing behaviors caused by cognitive overload in longer QSs. This research clarifies how human-centered design improves preference elicitation tools for collective decision-making.

## CCS Concepts

• **Human-centered computing** → **Collaborative and social computing systems and tools**; **Collaborative and social computing design and evaluation methods**; **HCI design and evaluation methods**; **Interactive systems and tools**; **Empirical studies in interaction design**.

## Keywords

Quadratic Survey; Preference Construction; Survey Response Format; Interactive User Interface; Cognitive Load

*Both authors contributed equally to this research.

## 1 Introduction

Designing intuitive survey interfaces is crucial for accurately capturing respondents' preferences, which directly impact the quality and reliability of the data collected. Recent Human-Computer Interaction (HCI) studies highlight how certain survey response formats can increase errors [43, 66] and influence survey effectiveness [93]. In this paper, our goal is to introduce an effective interface for a **Quadratic Survey (QS)**, a survey tool designed to elicit preferences more accurately than traditional methods [7]. Despite the promise of QSs, there has been no research on designing interfaces to support their unique quadratic mechanisms [31], where participants must rank and rate items — a task that poses significant cognitive challenges. To popularize QSs and ensure high-quality data, this paper addresses the question: *How can we design interfaces to support participants in completing Quadratic Surveys (QSs) more effectively?*

We envision an effective interface that navigates participants through the complex mechanism and preference construction process, tailored to a QS. A QS improves accuracy in individual preference elicitation compared to traditional methods like Likert scales by requiring participants to make trade-offs using a fixed budget of credits, where purchasing $k$ votes for an option in QS costs $k^2$ credits [7, 68]. This quadratic cost structure forces respondents to carefully evaluate their preferences, balancing the strength of their support or opposition against the limited budget. However, the process of making these thoughtful trade-offs introduces challenges. As individual preferences are often constructed when presented with the options [49], the act of weighing costs, evaluating options, and constructing rankings increases cognitive load. Moreover, QSs, often referred to as Quadratic Voting (QV) in voting scenarios, can

**Figure 1: The Two-Phase Interface: The interface consists of two phases. Survey respondents can navigate between phases using the top right button. In the organization phase, the interface presents one option at a time to the respondents, and they chose one of four positional choices: "Lean Positive", "Lean Neutral", "Lean Negative", or "Skip". Skipped options are hidden and can be evaluated later. The chosen options then appear below. Items can be dragged and dropped across categories or returned to the stack. In the voting phase, options are listed in the order of the four categories. When hovering over each option, respondents can select a vote for that option using a dropdown menu. Each dropdown menu contains the cost associated with the vote. A sort button allows ascending sorting within each category. A summary box tracks the remaining credit balance.**

involve hundreds of options [74, 87], increasing the risk of cognitive overload and the taking of mental shortcuts [63, 80, 92].

To date, existing quadratic mechanism-powered applications simply present options, allow vote adjustments and automatically calculate votes, costs, and budget usage. Such designs focused heavily on the mechanics operating the tool, rather than supporting possible challenges these application users faced. Survey interface literature, while addressing decision-making and usability, focuses on traditional surveys that do not share the unique option-to-option trade-offs that a QS introduces [19, 21, 41, 66, 91, 97]. Prior research in HCI and beyond explored techniques to manage cognitive load [50, 60, 62, 72, 91] and scaffold challenging tasks [36, 42, 46, 101] showing promise in supporting preference construction with a QS. Thus, this study aims to bridge this gap.

We propose a novel interactive two-phase "organize-then-vote" QS interface (referred to as the two-phase interface for short, Figure 1), which was developed through multiple iterations. It aims to facilitate preference construction and reduce cognitive load when making trade-offs through three key elements. First, the interface scaffolds the preference construction process by having participants initially categorize the survey options into "Lean Positive," "Lean Neutral," or "Lean Negative." This serves as a cognitive warm-up, easing participants into the more complex QS voting task. Second, the interface arranges the options according to these categorizations, providing a structured visual layout. Third, participants can refine the positions of these options using drag-and-drop functionality, giving them greater control and agency in the preference-construction process.

To explore how these interface elements affect cognitive load and support preference construction in QSs, we pose the following research questions:

**RQ1.** How does the number of options in Quadratic Surveys impact respondents' cognitive load?

**RQ2a.** How does the two-phase interface impact respondents' cognitive load compared to a single-phase text interface?

**RQ2b.** What are the similarities and differences in sources of cognitive load across the two interfaces?

**RQ3.** What are the differences in Quadratic Survey respondents' behaviors when coping with long lists of options across the two-phase interface and the single-phase text interface?

We invited 41 participants to a lab study comparing our two-phase interface with a baseline to understand how different interface designs and option lengths (6 options or 24 options) impact cognitive load.

Self-reported cognitive load using the NASA Task Load Index (NASA-TLX) and semi-structured interviews identified common challenges in QS, such as preference construction and budget management, while highlighting differences between text and two-phase interfaces. The two-phase interface fostered more strategic engagement with survey options, encouraging consideration of broader impacts in the long QS, reducing time pressure in the short QS, and eliciting greater affirmative satisfaction (e.g., "feeling good"). Quantitative results support these observations: participants in the two-phase interface—particularly in long surveys—traversed the list less frequently but maintained the same number of edits while spending more time per option. This suggests that reduced traversal

did not diminish engagement. Together, these findings highlight the organizing phase's role in fostering deeper engagement with survey options.

*Contributions.* We contribute to the body of knowledge in the HCI community by proposing the first interface specifically designed for QS and QV-like applications, which aims to reducing cognitive challenges and scaffolding preference construction through a two-phase interface with direct manipulation. Before our work, no research had explored QS interfaces. This is particularly important for long QSs, which are prone to cognitive overload. Few HCI studies have addressed interfaces for surveys and questionnaires. Our study demonstrates how user interfaces can facilitate preference construction in situ and promote deeper engagement with survey options through interface elements. Additionally, this paper offers the first in-depth qualitative analysis of user experiences with Quadratic Mechanism applications, identifying usability challenges and key factors contributing to cognitive load. The impact of our contribution extends beyond QSs, offering design implications for other preference-elicitation tools used in multi-option scenarios. By making QSs easier to use and more accurate, our design encourages wider adoption among researchers and practitioners. Finally, our work lays the groundwork for future Quadratic Mechanism interface design to facilitate individuals expressing their preferences.

## 2 Related Work

This research lies at the intersection of three core areas: quadratic surveys, existing QV interfaces and choice overload along with cognitive challenges. In this section, we review the related works in each of these areas.

## 2.1 Quadratic Survey and the Quadratic Mechanism

We introduce the term **Quadratic Survey (QS)** to describe surveys that utilize the quadratic mechanism to collect individual attitudes. The **quadratic mechanism** is a theoretical framework designed to encourage the truthful revelation of individual preferences through a quadratic cost function [31]. This framework gained popularity through **Quadratic Voting (QV)**, also known as plural voting, which uses a quadratic cost function in a voting framework to facilitate collective decision-making [47].

To illustrate how QS works, we formally define the mechanism: each survey respondent is allocated a fixed budget, denoted by $B$, to distribute among various options. Participants can cast $n$ votes for or against option $k$. The cost $c_k$ for each option $k$ is derived as:

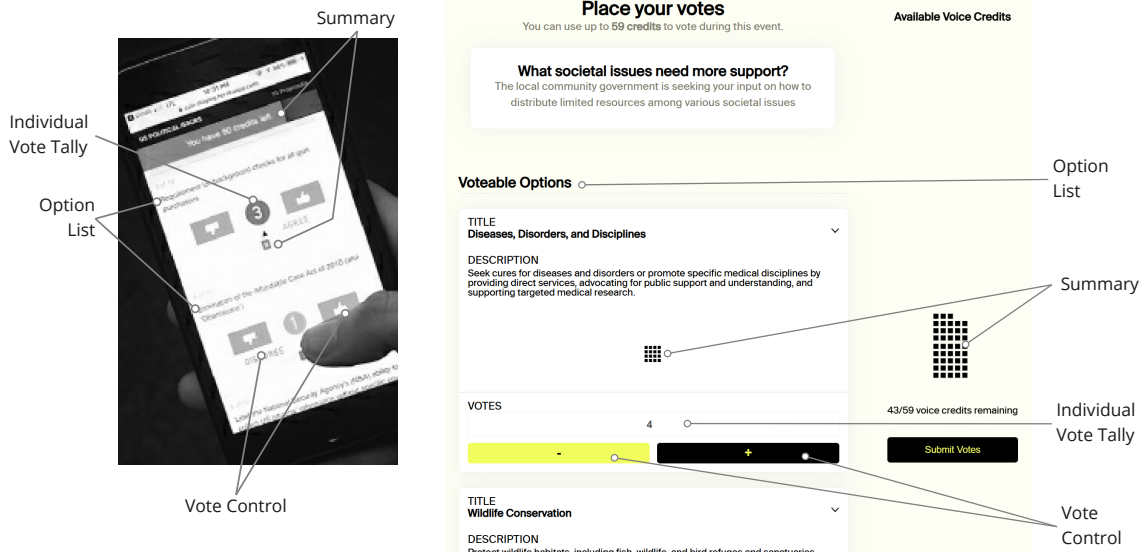$$c_k = n_k^2 \quad \text{where} \quad n_k \in \mathbb{Z}$$

The cost of all votes must not exceed the participant's budget:

$$\sum_k c_k \leq B$$

Survey results are determined by summing votes for each option:

$$\text{Total Votes for Option } k = \sum_{i=1}^{S} n_{i,k}$$

where $S$ represents the total number of participants, and $n_{i,k}$ is the number of votes cast by participant $i$ for option $k$. Each additional

**Figure 2: A selection of two QV interfaces. The interface on the left was used in the first empirical QV research [68]. Little information is available about the software, except for an image from Posner and Weyl [67]. The interface on the right is an open-sourced QV interface [71] forked from GitCoin [30], used by the RadicalxChange community [70]. Both interfaces share the common elements with different visual representations.**

vote for each option increases the marginal cost linearly, encouraging participants to vote proportionally to their level of concern for an issue [67].

QS adapts these strengths of the quadratic mechanism in *voting* to encourage truthful expression of preferences in *surveys*. Unlike traditional surveys that elicit either rankings *or* ratings, QS allows for *both*, enabling participants to cast multiple votes for or against options, incurring a quadratic cost. Cheng et al. [7] showed that this mechanism aligns individual preferences with behaviors more accurately than Likert Scale surveys, particularly in resource-constrained scenarios like prioritizing user feedback on user experiences.

In recent years, empirical studies on QV have expanded into various domains [4, 56]. Applications based on the quadratic mechanism have also grown, including Quadratic Funding, which redistributes funds based on outcomes from consensus made using the quadratic mechanism [2, 26]. Recent work by South et al. [82] applies the quadratic mechanism to networked authority management, later used in Gov4git [18]. Despite the increasing breadth and depth of applications utilizing the quadratic mechanism, little attention has been paid to user experience and interface design, which support individuals' preference intensity elicitation. Our work aims to address this by designing interfaces for quadratic mechanisms.

## 2.2 Existing QV Interfaces

Since QS shares QV's underlying mechanism, we used snowball sampling to identify publicly available QV applications mentioned in news and academic sources. Currently, no widely adopted QV interface is tied to a single vendor or platform. Figure 2 shows two variations of existing QV interfaces, with both employing a single-step approach with different visual representations of common elements [4, 7, 18, 100]. All QV interfaces generally include:

- Option list: A list of options for voting.
- Vote controls: Buttons to add or remove votes for each option.
- Individual vote tally: A display of the votes cast per option.
- Summary: A summary of costs and the remaining budget.

These components let users operate QV mechanically, providing little understanding of voters' usability needs nor offering cognitive support. In addition, HCI research on survey interfaces is limited [57, 94] with most efforts focusing on alternative input modalities like bots, voice interfaces, or virtual reality [23, 40, 43, 96].

## 2.3 Cognitive Challenges and Choice Overload

The challenge of respondents making difficult decisions using quadratic mechanisms remains unexplored in the literature. Lichtenstein and Slovic [49] identified three key elements that make decisions difficult. These elements include making decisions in unfamiliar contexts, quantifying the value of one's opinions, and being forced to make trade-offs due to conflicting choices. QS fits at least two of the three elements: participants may encounter a selection of unfamiliar options by the survey designer; they are asked to quantify the difference between option preferences through a numerical vote; and the budget constraint enforces trade-offs under a non-linear function, which means that a vote decrease for one option is not necessary equivalent to an increase for another, making iterative adjustment and evaluating tradeoffs difficult. Thus, we believe QS introduces a high cognitive load.

Cognitive load refers to the demands placed on a user's working memory during the interaction process, which significantly influences the usability of the system [12, 78]. Cognitive overload can adversely affect performance [17], leading individuals to rely on heuristics rather than deliberate, logical decision-making [15]. When presented with excessive information, such as too many options, individuals 'satisfice', settling for a 'good enough' solution rather than an optimal one [63, 80, 92]. Subsequently, too many options can overwhelm individuals, resulting in decision paralysis, demotivation, and dissatisfaction [37].

Additionally, Alwin and Krosnick [1] highlighted that the use of ranking techniques in surveys can be time-consuming and potentially more costly to administer. These challenges are compounded when ranking numerous items, requiring substantial cognitive sophistication and concentration from survey respondents [22].

Notable applications of QV include the 2019 Colorado House, which considered 107 bills [14], and the 2019 Taiwan Presidential Hackathon, which featured 136 proposals [65]; both used a single QV question with hundreds of options. These empirical applications of QV suggest the importance of understanding QS with many options' impact on cognitive load and support developing interfaces for practical uses.

## 3 Quadratic Survey Interface Design

This section presents our QS interface. Drawing on existing QV interfaces described in Section 2.2 and prior literature, we iterated through paper prototypes and three design pre-tests, detailed in Appendix A. Initially, participants struggled to *rank* relative preferences among options and *rate* the degree of trade-offs between them. In this study, we focus on addressing the former challenge, which pertains to preference construction.

### 3.1 'Organize-then-Vote': The Two-Phase Interface

*3.1.1 Justifying a two-phase approach.* The main objective of the two-phase interface is to facilitate preference construction and reduce cognitive load. As shown in Figure 1, the interface consists of two steps: an organization phase and a voting phase. In both phases, survey respondents can drag and drop options across the presented list.

*A two-phase approach.* Preferences are shaped through a series of decision-making processes [49]. Two decision-making theories inspired this two-step interaction interface design: Montgomery [53]'s Search for a Dominance Structure Theory (Dominance Theory) and Svenson [85]'s Differentiation and Consolidation Theory (Diff-Con Theory). The former suggested that decision-makers prioritize creating dominant choices to minimize cognitive effort by focusing on evidently superior options [53]. The latter described a two-phase process where decisions are formed by initially *differentiating* among alternatives and then *consolidating* these distinctions to form a stable preference [85]. Pre-tests showed participants puzzled by ranking all options before voting. These theories suggest decisions emerge by eliminating choices, not by fully ranking them. Therefore, the organize-then-vote design makes this natural process more explicit. Phase one focused on differentiating and identifying dominant options, enabling survey respondents

to preliminarily categorize and prioritize their choices. Phase two presented these categorized options in a comparable manner, with drag-and-drop functionality, enhancing one's ability to consolidate preferences. This structured approach aimed to construct a clear decision-making procedure that reduced cognitive load and enhanced clarity and confidence in the decisions made.

*Phase 1: Organization Phase.* The goal of the organization phase was to support participants in identifying clearly superior options or partitioning choices into distinguishable groups. In this section, we first describe how the interaction works, then we detail the reasons for the implemented design decisions.

The organizing interface, depicted on the top half of Figure 1, sequentially presents each survey option. Participants select a response among three ordinal categories – "Lean Positive", "Lean Negative", or "Lean Neutral". Once selected, the system moves that option to the respective category. Participants can skip the option if they do not want to indicate a preference. Options within the groups are draggable and rearrangeable to other groups should the participants wish.

To support preference formation, respondents are shown one option at a time, allowing them to either recall a prior judgment or construct a new one based on the presented choices [83]. Limiting the information presented this way also helps reduce cognitive load by preventing overload from too many options [86]. This incremental process ensures that participants form opinions on individual options.

The three possible options — Lean Positive, Lean Neutral, and Lean Negative — aim to scaffold participants in constructing their own choice architecture [55, 88], which strategically segments options into diverse and alternative choice presentations while avoiding biases from defaults. We believed that these three categories were sufficient for participants to segment the options. We do not limit the number of options one can place in each category to prioritize user agency, allowing participants full control over how they organize their preferences [58]. Immediate feedback displays the placement of options and allows participants to rearrange them via drag-and-drop, adhering to key interface design principles [58]. It also allows finer-grain control for individuals to surface dominating options and create differentiating groups of options.

*Phase 2: Interactive Voting Phase.* The objective of the voting phase is to facilitate the consolidation of differentiated options through interactive elements while reinforcing the differentiation across options constructed by participants in the previous phase. This facilitation is achieved by retaining the drag-and-drop functionality for direct manipulation of position and enabling sorting within each category.

Options are displayed as they are categorized within each category from the previous step and in the following section — Lean Positive, Lean Neutral, Lean Negative, and Skipped or Undecided — as detailed on the bottom half of Figure 1. The Skipped or Undecided category contains options left in the organization queue, possibly because survey respondents have a pre-existing preference or chose not to organize their thoughts further. The original order within these categories is preserved to maintain and reinforce the differentiated options. This ordering sequence mitigated early prototype concerns where uncategorized options were left at the top

of the voting interface confusing survey respondents. Respondents have the flexibility to return to the organization interface at any point during the survey to revise their choices.

In the voting interface, options are draggable, allowing participants to modify or reinforce their preference decisions as needed. Each category features a sort-by-vote function for reordering within the group, which, although it doesn't affect the final outcome, supports information organization and consolidation. Both features aim to group similar options automatically and emphasize proximity, reducing cognitive load by following the proximity compatibility principle to enhance decision-making [99].

While multiple interaction mechanisms exist, drag-and-drop has been extensively explored in rank-based surveys. For instance, Krosnick et al. [44] demonstrated that replacing drag-and-drop with traditional number-filling rank-based questions improved participants' satisfaction with little trade-off in their time. Similarly, Timbrook [89] found that integrating drag-and-drop into the ranking process, despite potentially reducing outcome stability, was justified by the increased satisfaction and ease of use reported by respondents. The trade-off was deemed worthwhile as QSs did not use the final position of options as part of the outcome if it significantly enhanced user satisfaction and usability [73]. Together, these design decisions led to our belief that a two-phase interface with direct interface manipulation could reduce the cognitive load for survey respondents to form preference decisions when completing QSs.



**Figure 3: Alternative vote control. The click-based design (upper) mirrors traditional vote control used in other QV interfaces, where each click controls one vote. The wheel-based design (the latter two) allows control through both clicks and mouse wheel rotation.**

In addition, we made three aesthetic design decisions considering existing QV-based interfaces. First, we removed visual elements like icons, emojis, progress bars, and vote visualizations, as prior research indicated that emojis could influence survey interpretations and reduce user satisfaction [35, 91]. While effective visualizations can aid decision-making, this study does not aim to address that question. Second, all options are visible on the screen simultaneously. Prior research recommends placing all items on the voting screen to prevent overlooked votes [5]. This echoes the proverb "out of sight, out of mind," reducing where individuals might be biased toward visible options, and additional effort is required for individuals to retrieve specific information if options are hidden. Last, use a dropdown positioned to the right of each survey option

for ease of access to the budget summary when determining the votes. The layout of the votes and cost was inspired by online shopping cart checkout interfaces where quantities are supplied next to the itemized costs followed by the total checkout amount. Figure 3 shows the two alternatives—click-based buttons (participants disliked multiple clicks) and a wheel-based design (unfamiliar to some)—and settled on the dropdown.

## 3.2 Baseline Interface: Single-Phase Text Interface

We created a single-phase text interface (referred to text interface for short, Figure 4) as a control, enabling us to see how organizational features affect cognitive load and behavior. Like existing interfaces, it uses static lists, a summary box, and a vote control. To ensure a fair comparison, we applied the same design principles: no extraneous visuals, all options on one screen, and dropdown-based voting. The prompt appears at the top, followed by a randomly ordered list to prevent ordering bias [13, 24]. Costs and the credits summary appear on the right.

Both experimental interfaces were developed with a ReactJS frontend and a NextJS backend powered by MongoDB. We open-source both interfaces.[1]

## 4 Experiment Design

In this section, we describe our experiment design. The study was approved by the university's Institutional Review Board (IRB).

### 4.1 Recruitment and Participants

We recruited 41 participants from a United States college town using online ads, digital bulletins, social media posts, email newsletters, and physical flyers in public spaces beyond campus. We described the study as exploring societal attitudes to reduce response bias. One participant was excluded due to data quality concerns[2].

To ensure diversity, we prioritized non-students by selectively accepting them and monitoring demographic distribution. The mean participant age was 34.63 years, with an age distribution similar to the county's demographic profile (Figure 5a), although there was a slightly higher representation of younger adults. Gender and race demographics are presented in Figures 5b and 5c. Demographic differences between groups were reasonably balanced, although participants using the short text interface skewed slightly younger ($\mu$=32.1), and those in the long two-phase interface group had a broader age range ($\mu$=38.8, $\sigma$=19.6). Appendix D contains full details.

### 4.2 Experiment Design

We implemented a between-subject design to avoid learning effects and minimize participants' fatigue from potential complexity of QSs. The experiment focused on public resource allotment, following the methodology of Cheng et al. [7], in which participants expressed preferences across societal issues. These issues are relevant to all citizens and effectively highlight the need to prioritize limited public resources. Participants received a survey with options randomly

---

## What societal issues need more support?

The local community government is seeking your input on how to distribute limited resources among various societal issues. Using the **quadratic survey mechanism**, please indicate your preferences below. *Upvote more* for issues you think deserve more resources, and *downvote more* for those you believe should receive fewer resources.

You have 59 credits to distribute. You can vote on each option by clicking the dropdown menu when you hover over the option.

**All Options**

| | | |
|---|---|---|
| **Youth Education Programs and Services** Provide programming, classroom instruction, and support for school-aged students in various disciplines such as art education, STEM, outward bound learning experiences, and other programs that enhance formal education. | No votes | $0 |
| **Advocacy and Education** Support social justice through legal advocacy, social action, and supporting laws and measures that promote reform and protect civil rights, including election reform and tolerance among diverse groups. | 3 upvotes | $9 |
| **Zoos and Aquariums** Support and invest in zoos, aquariums and zoological societies in communities throughout the country. | 6 upvotes | $36 |
| **Community Foundations** Promote giving by managing long-term donor-advised charitable funds for individual givers and distributing those funds to community-based charities over time. | 2 downvotes | $4 |
| **Environmental Protection and Conservation** Develop strategies to combat pollution, promote conservation and sustainable management of land, water, and energy resources, protect land, and improve the efficiency of energy and waste material usage. | 1 upvote  $1  ⌄ | |
| **International Peace, Security, and Affairs** Promote peace and security, cultural and student exchange programs, improve relations between particular countries, provide foreign policy research and advocacy, and United Nations-related organizations. | 5 upvotes  $25 4 upvotes  $16 3 upvotes  $9 2 upvotes  $4 **1 upvote  $1** No votes  $0 | |

Non-draggable Randomly Positioned Options

Hover and click to show vote options

Budget Summary

**Credit Summary**

**Remaining Credit**  $9

Submit

**Figure 4: The text-based interface: This interface is based on the two-phase version but does not include the organization phase and lacks the drag-and-drop functionality. Options are randomly positioned.**

drawn from the 31 societal topics evaluated by Charity Navigator [6], an organization that assesses over 20,000 charities in the United States (see Appendix C for the full list). Randomly selecting the options each participant saw helped control for potential systematic content biases that specific voting options might introduce across surveys of different lengths. Participants were randomly assigned to one of four groups, each with 10 participants:

- Short Text (ST): A text interface with 6 options.
- Short Two-Phase (S2P): A two-phase interface 6 options.
- Long Text (LT): A text-based interface 24 options.
- Long Two-Phase (L2P): A two-phase interface with 24 options.

Prior research informed the choice of 6 and 24 options, representing short and long lists. These studies recommend fewer than 10 options for constant-sum surveys [54] and fewer than 7 for the Analytic Hierarchy Process [76]. Classic cognitive load research [52, 77] suggests the use of 7±2 items. A meta-analysis by Chernev et al. [8] identified 6 and 24 as common values for short and long lists in choice overload studies, which are rooted in the original choice overload experiment by Iyengar and Lepper [37].

## 4.3 Experiment Procedure

Participant's spent on average 40 minutes (range: $27-68$, $\sigma$=9) in the lab. Figure 6 visually represents the study protocol detailed in the following subsections.

*4.3.1 Consent, Instructions, and Quiz.* Participants were invited to the lab to control for external influences and used a 32-inch vertical monitor to display all options. After consenting, participants watched a video explaining the quadratic mechanism without any mention of the interface's operation, followed by a quiz to ensure understanding. Participants rewatched the video or consulted the researcher until they successfully selected the correct answers. Each participant's screen was captured throughout the study.

*4.3.2 Quadratic Survey.* The researcher informed participants that the study aimed to help local community organizers understand preferences on societal issues to improve resource allocation. Aware that their screens were being recorded, participants completed the survey independently inside a semi-enclosed space in the lab, without the researcher's presence. Once they completed the survey, participants notified the researcher.

*4.3.3 NASA-TLX Survey and Interview.* The researcher joins study participant and administer a paper-based weighted NASA Task Load Index (NASA-TLX), followed by a semi-structured interview after being informed that the researcher would begin audio recording with their laptop. We adopted the paper-based weighted NASA-TLX, a widely used multidimensional tool that averages six sub-scale scores to measure overall workload after task completion [3,

(a) Age distribution of the study participants were similar to the locale's demographic profile.



(b) Gender distribution of our participants skewed towards female participants.



(c) Ethnicity distribution remains diverse with fewer Hispanic and African American participants.

Figure 5: Demographic distributions: Age, Gender, and Ethnicity



Figure 6: Study protocol: Participants are asked to learn about the mechanism of QSs after consenting to the study. The researcher explained the study overview and asked participants to complete the QS. A NASA-TLX survey followed by interviews to understand participants' cognitive load. We debriefed participants after the study.

33, 34]. NASA-TLX is favored for its low cost and ease of administration [27], and it exhibits less variability compared to one-dimensional workload scores [75], making it suitable for our study.

While cognitive load can be assessed through psychophysiological, performance, subjective, and analytical measures [27], the length and complexity of QSs make some of these impractical. Performance and analytical measures require task switching or interruptions, which risk increasing overall cognitive load and experiment time. Psychophysiological measures, such as pupil size [61]

and ECG [32], are costly, sensitive to external factors, and often require participants to wear additional equipment.

*4.3.4 Demographic, Debrief, and Compensation.* After the interview, the researcher collected participant's demographics and debriefed them, explaining that the study's goal was to understand interface design and cognitive load. Participants received a $15 cash compensation.

## 4.4 Quantitative Measures: Clickstream Data

Besides using NASA-TLX and interviews to capture cognitive load, we also recorded participants' clickstream data from the interface (i.e., each click and the corresponding UI component). These log data enabled us to analyze how participants navigated and engaged with the survey options.

*Edit Distance.* We introduce three related metrics—edit distance per option, edit distance per action, and cumulative edit distance—to quantify the distance participants traveled across the interface. Edit distance per option sums the total number of options traversed while modifying a single vote option, edit distance per action measures the distance traversed during each individual adjustment, and cumulative edit distance captures the total distance traversed throughout the entire survey. The formal definitions and modeling approach are provided in Section 6.

*Time Spent per Option.* In addition, we computed the total time participants spent interacting with each specific option by aggregating the time spent on that specific option during the survey. We describe and discuss these findings in Section 7.

## 5 Result: Self-Reported Cognitive Load in Quadratic Surveys

This section presents findings on cognitive load in QSs, focusing on how the number of options and different interfaces influence it (**RQ1**, **RQ2a**). We analyze similarities and differences in cognitive load sources across conditions (**RQ2b**).

Qualitative findings are based on an inductive thematic analysis [59], which was conducted after transcribing the interviews. The first author single-coded the snippets according to the research questions and merged them into overarching themes. The first author conducted multiple rounds of coding, and identified differences across conditions, which were refined and validated using a deductive coding process.

Quantitative findings are derived from a Bayesian approach, which enhances transparency by interpreting posterior distributions and moving beyond binary thresholds [39]. Bayesian methods suit various sample sizes, leveraging maximum entropy priors to ensure conservative and robust inferences [51].

## 5.1 Overall Cognitive Load from NASA-TLX

Weighted NASA-TLX uses a continuous 0 to 100 score, with higher values denoting greater cognitive load. We use predefined mappings of NASA-TLX scores to cognitive levels: low, medium, somewhat high, high, and very high, as described by Hart and Staveland [34]. Results are shown in Figure 7a, with value interpretations presented in Figure 7b.

Given the sparsity of the data, we modeled the weighted NASA-TLX scores as ordinal outcomes based on value interpretations. We developed a hierarchical Bayesian ordinal regression model with length as an ordinal predictor and interface type as a categorical predictor, using hierarchical priors for partial pooling. Interaction effects between length and interface are captured via a non-centered parameterization with an LKJ prior to account for correlations [51]. We applied the same model to the NASA-TLX subscales; since these subscales lack inherent cognitive level interpretations, we

constructed weighted bins for the ordinal regression. In our model, a latent variable represents a continuous measure of cognitive load, discretized into ordinal outcomes via thresholds. Details of this model and additional subscale results are provided in Appendix G.

In Bayesian analysis, the 94% high-density interval (HDI) represents the range where the true parameter is most likely to lie. While the results (Figure 8) in terms of differences in latent cognitive load are not statistically significant because 0 is within this range, the HDI quantifies probabilistic trends and accounts for uncertainty in a transparent manner.

- Increased option length with text interface trends to *reduced* cognitive load with a posterior probability of approximately 84.5%. This reflects a median cognitive load of 33.85 (mean = 34.60, SD = 17.69) compared to a median of 39.00 (mean = 43.23, SD = 17.65).
- Within short QSs, the two-phase interface trends to *reduced* cognitive load, with a posterior probability of 77.6% supporting the reduction. Participants report a median cognitive load of 29.85 (mean = 35.36, SD = 18.17) under the two-phase interface compared to a median of 39.00 (mean = 43.23, SD = 17.65) under the text interface.
- For the long QSs, the two-phase interface trends an *increase* in cognitive load with a posterior probability of 62.7%. The median cognitive load is 42.70 (mean = 42.02, SD = 18.48) under the two-phase interface compared to 33.85 (mean = 34.60, SD = 17.69) in the text interface.

This result contradicts our hypothesis that more options would increase cognitive load and that interfaces can reduce it. Thus, we explore qualitative results to identify possible explanations. To understand the similarities and differences in sources of cognitive load (**RQ2b**), we analyze qualitative results across the six NASA-TLX subscales: mental demand, physical demand, temporal demand, effort, frustration, and performance. Detailed breakdown of each subscale are provided in Appendix E.

## 5.2 Qualitative Analysis: Common Sources of Cognitive Load

Our analysis reveals several themes across different cognitive load subscales. We focus on three themes common to all experimental conditions, omitting less related themes for clarity.

**Preference Construction** is cited by 97.5% (N=39) of participants as a significant source of mental demand, consistent with prior literature suggesting that preferences are often constructed in context rather than fixed [49]. Specific tasks contributing to this demand include evaluating the relative importance between options (e.g., S002 💬 *Figuring out[…] how much I prioritize option 1 over option 2* , 40% (N=16)), making trade-offs due to limited resources (e.g., S005 💬 *[…] very hard to take decisions …I felt that multiple options deserve equal amounts of credit …but you have given very limited credit.* , 42.5% (N=17)), and deciding the exact number of votes (e.g., S023 💬 *[…] having to pick how many upvotes would go to each one* , 70% (N=30)).

**Budget Management** emerges as a source of both mental and temporal demand. 25% (N=10) of participants describe the challenge of working with limited credits while trying to maximize their allocation (e.g., S032 💬 *[…] for certain societal issues, you had to …take away from other issues you could support* ). An equal

(a) NASA-TLX Weight Score: The Long Two-Phase Interface exhibits the highest weighted cognitive load with a median of 42.70, a mean of 42.02. This is higher than the long text interface, which has a median cognitive load of 33.85 and a mean of 34.60. However, the short text interface demonstrates a higher cognitive load with a median of 39.00, a mean of 43.23, compared to the short two-phase interface, which has a median of 29.85, a mean of 35.36. The standard deviation is similar across groups at around 18.

(b) NASA-TLX Cognitive Interpretation: More participants in the short text interface, totaling 8, reported a somewhat high or above cognitive load, which is significantly higher compared to the 5 participants who reported similarly for the short two-phase interface. However, the long two-phase interface saw slightly more participants, 6 in total, reporting somewhat high or above cognitive load compared to the long text interface.

**Figure 7: This figure shows the box plot results for weighted NASA-TLX scores across experiment groups and participant counts based on individual score interpretations. In 7a, we observe a downward trend in cognitive load for the short QS, while the long QS shows an upward trend. Interestingly, there is a counterintuitive downward trend between short and long text interfaces. In 7b, these trends are clearer when NASA-TLX scores are grouped into five tiers.**



**Figure 8: Posterior distributions of differences in latent cognitive load between experimental conditions. Values below 0 indicate reduced load. Main takeaway: while the model does not indicate statistically significant differences, longer text interfaces are more likely to reduce cognitive load, and the two-phase interface has a higher probability of lowering cognitive load.**

percentage of participants find it mentally taxing to keep track of remaining credits (e.g., S006 ✎ *[…] looking at the remaining credits, I'm trying to mentally divide that up before I start allocating* ).

When assessing themes across all subscales, we identified patterns that highlights the underlying nature of participants' cognitive efforts across different contexts. Thus, we also coded interview snippets as **Operational** and **Strategic** actions in addition to

goal-oriented actions such as Budget Management and Preference Construction.

**Operational Actions** refer to reactive efforts addressing immediate, tactical needs, which emerged across all experimental conditions. These actions involve direct task execution, responding to constraints without reflection on broader, long-term implications. Examples include adjusting choices to stay within budget

(e.g., S003 💬 *I had to alter [...] I kept going under budget* ), re-reading options (e.g., S010 💬 *I just had to reread it again* ), completing questions efficiently (e.g., S010 💬 *I was trying to be efficient in responding to the question* ), and interacting with the survey interface (e.g., S018 💬 *Like (deciding) one upvote or two upvotes[...]* ). 40% (N=16) of participants attribute Operational actions to temporal demand. Additionally, 37.5% (N=15) attribute this cause to frustration, and 32.5% (N=13) attribute it to performance. While commonly cited across conditions, its distribution varies.

## 5.3 Qualitative Analysis: Different Sources of Cognitive Load

There are several notable differences between the text and two-phase interfaces.

First, regardless of length, when analyzing performance, which refers to a person's perception of their success in completing a task, participants describe their performances differently. We categorize them into indications of satisficing behaviors("good enough"), exhausting their effort (i.e., "done their best,"), or feeling positive (i.e., "feeling good.") There are almost twice as many participants using the two-phase interface to report a positive feeling about their final submission (55% v.s 30% (N=11 vs. 6)).

Second, 70% (N=14) of text interface participants attribute operational actions as contributors to effort, double the percentage observed in the two-phase interface group (35%, N=7). This partially echoes the finding that 90% (N=18) of text interface participants report mental demand from deciding the exact number of votes, compared to 60% (N=12) in the two-phase interface group.

The distinction between the text and two-phase interfaces becomes more pronounced in the context of the long survey. 80% of the long text interface participants (N=8) attribute operational actions to effort, compared to only 20% (N=2) in the long two-phase interfaces. Conversely, 90% of long two-phase interface participants (N=8) attribute effort to strategic actions, compared to 50% (N=5) in the text interface.

We also found differences in how preference construction differs in contributing to their mental demand and sources of effort. Opposite to operational actions, **strategic considerations** refer to considering about long term goals, determining priorities, considering broader implications, and considering option's more holistically. Consider the following quotes:

> *Trying to figure out what upvotes I should give [...] went back and forth between those two. [...] it was very mentally tasking for me.*
> 💬 S015 (LT)

> *[...] especially with so many different societal issues. How do I personally prioritize them? And to what extent do I prioritize them?*
> 💬 S009 (L2P)

S015 describes the **operation** of locating tasks to find the right vote, whereas S009 **strategically** aligns higher-order values holistically. Regarding mental demand, 80% of participants in the long text interface focused on a narrower scope, comparing fewer options (N=8), while only 30% did so in the two-phase interface (N=3). Conversely, 90% of participants in the long two-phase interface considered broader societal impacts and evaluated more options simultaneously (N=9), compared to 30% in the text interface (N=3). Similar distinctions were evident in effort-related sources.

These differences highlight variations in **levels of engagement** with the survey content. Participants using the two-phase interface expressed higher satisfaction with their performance. For the long survey, they engaged with broader aspects across different options and strategically allocated their credits.

## 5.4 Qualitative Analysis: Instances of Satisficing

When individuals cannot process all available information, prior research has found that people exhibit *satisficing behaviors*, which refers to settling for *good enough* rather than *optimal* decisions [28]. While we did not explicitly ask participants if they 'satisfied,' nor did we measure it quantitatively, we identified satisficing behaviors based on participants' explanations of how they completed the survey. For example,

> *[...] you thought of enough things, you know, and so it wasn't the most effort I could put in because again, that would have been diminishing returns. I tried to think of enough things [...] and then move on. [...]*
> 💬 S032 (ST)

> *I felt like that (the response) was satisfied, but not perfect. Cause perfect is not a reality.*
> 💬 S036 (ST)

This quote illustrates satisficing decision-making, where participants chose to settle for suboptimal outcomes. Satisficing was observed primarily at the beginning and end of the survey, where participants allocated large amounts of credit initially and then managed the remaining credits to confirm their final vote allocations. For instance,

> *[...] Because that (the credit) was what was left. [Laughter] I probably wouldn't use that on <optionA> instead of <optionB>. [...]*
> 💬 S015 (LT)

> *[...] it went negative, and then I just settled for just $6 remaining. [...] I don't think it's perfect. But I think I'm satisfied. Yeah, I'm satisfied.*
> 💬 S033 (LT)

> *[...] when I had first started like looking at the first few, I was just doing it kinda like willy nilly, I'm not really paying that much attention to necessarily how many credits I had, or how many categories there were.*
> 💬 S041 (LT)

Participants also exhibited satisficing behaviors regarding *defaults*, particularly when constructing their preferences. For example, participant S003, described how default placements influenced their final decisions:

> *Honestly, if medical research [...] was the first one I saw, I think it would automatically give it a lot more.*
> 💬 S003 (ST)

Our qualitative analysis found that 60% of short-text participants (N=6) and 50% of long-text participants (N=5) expressed instances of satisficing behaviors when describing how they completed the survey, compared to none of the short two-phase participants and 30% of long-text participants (N=3). These qualitative results highlighted potential satisficing behaviors across conditions.

## 6 Clickstream data: Interface reduces edit distance in long surveys

Following our findings on cognitive load, we analyze voting behaviors to identify differences in how participants cope with survey lengths, how interfaces influence their behavior, and why the long

text interface might exhibit lower cognitive load. All data are publicly available[3] to ensure transparency and support further research. This measure reveals how participants navigate and engage with survey options. We examine three dimensions of this measure: edit distance per option, edit distance per action, and cumulative edit distance throughout the survey.



**Figure 9: Edit Distance Per Option: We sum the total number of edit distances for each option, with the figure using the radius to indicate how often a specific edit distance occurred within an experimental condition. Main takeaway: Participants in the two-phase interface completed their votes for more options with fewer edit distances, whereas the Long Text interface shows a long tail of options requiring a wider range of edit distances.**

**Edit distance per option:** We calculate the total number of options a participant traversed when adjusting votes for a single option. Figure 9 illustrates differences across experimental conditions, with the long text interface showing the largest variance in the distance traveled and the highest mean. We implement a hierarchical Bayesian framework to model edit distance differences across experimental conditions. The observed distance differences are modeled using an exponential distribution, where the scale parameter is linked to survey length (treated as an ordinal variable), interface type (treated as a categorical variable), interaction effects between length and interface, and controlling for individual user variability. The linear predictor includes a global intercept and slope for length, random effects for each interface condition with an LKJ prior that captures the correlations among interface categories, and user-specific random effects to account for individual heterogeneity. Appendix I.1 includes the detailed model.

Figure 10 illustrates the pairwise posterior distributions for differences in edit distances across experimental conditions. For example, the difference in edit distances between the short and long static interfaces has a mode of 9.1, with a 94% highest density interval (HDI) of [6, 13]. This indicates that participants in the long text interface move approximately 9.1 steps more than those in the short text interface, with a high degree of confidence. The effect size is large (mode = 5.1, 94% HDI = [3.3, 7.1]), suggesting a statistically

[3]https://github.com/CrowdDynamicsLab/Quadratic-Survey-Dataset-and-Analysis

significant difference, which is expected due to the greater number of options in the long text interface.

Similarly, two-phase interface participants make approximately 8.9 fewer steps per option (mode = 8.9, 94% HDI = [6.4, 12]) than those in the long text interface, with a large effect size (mode = 5.7, 94% HDI = [4.2, 7.9]). The increase in edit distances between the short and long two-phase interfaces is substantially smaller (mode = 1.7, 94% HDI = [-0.01, 3.1]) compared to their static counterparts. Comparing the short text and short two-phase interfaces shows limited difference (mode = 1.3, 94% HDI = [-0.78, 3.8]), though the posterior distribution favors fewer steps for the two-phase interface (89.3% probability). The model suggests that the two-phase interface reduces edit distance per option, particularly for the long QS.
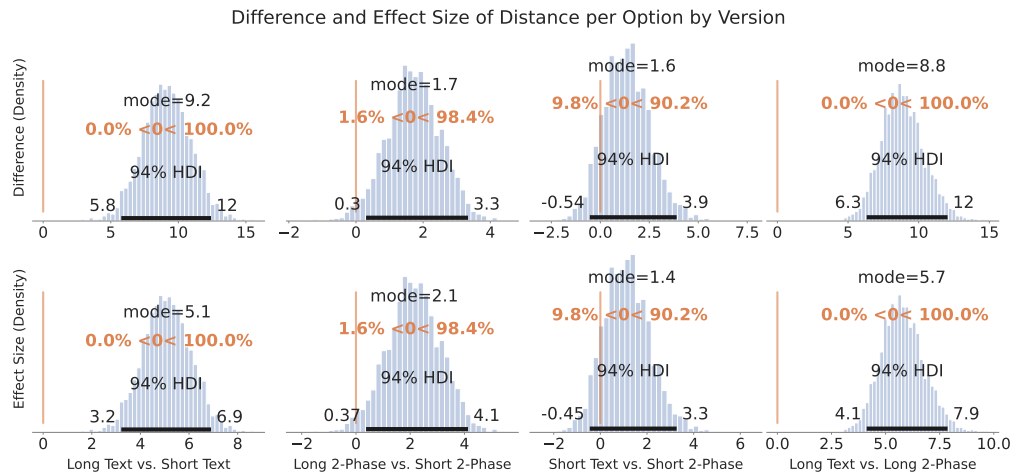
**Edit distance per action:** Building on the statistical disparities observed in the previous analysis and the unique patterns exhibited by long text interface participants, we present analyses focusing on edit distance per action and cumulative edit distance throughout the survey between the long text and long two-phase interfaces. Edit distance per action measures how far participants move during each adjustment while completing the survey. Figure 11 illustrates how, at each step, the number of participants moving a given distance (represented by the size of the dots) varies across experimental conditions. Visually, participants move less on average per option within the two-phase interface, with lower variance at smaller scales. This indicates that participants are making local edits, meaning their adjustments tend to occur near their previous edits in terms of edit distance. This also highlights that the organization phase effectively adjusts option positions for easier access, despite participants still having the freedom to move across the interface as all options are presented to them.

In contrast to earlier analyses, we use a hierarchical Bayesian model (detailed in Appendix I.2) to jointly estimate the mean and variance of edit distances across experimental conditions. The model assumes that edit distances are continuous and follow a normal distribution. This approach accounts for both central tendencies and variability, using separate predictors for the mean and variance. The model includes hierarchical effects for survey length, interface type, interactions between length and interface, and user-level random effects. Non-centered parametrization is used for survey length and interface type to improve convergence, while interaction effects are modeled with an LKJ prior to capture the correlations between factors.

Figure 12 illustrates the posterior variance distributions, confirming our hypothesis. Participants in the long text interface exhibit greater variance in movement, frequently navigating across the interface, compared to those in the long two-phase interface. This is evidenced by a variance difference mode of 76 (95% HDI = [59, 99]) and a large effect size (mode = 7.1, 95% HDI = [5.5, 9.2]).

**Cumulative edit distance for a participant:** Figure 13 illustrates how the two-phase interface reduces per-action distance, accumulating over time. Some long text participants traverse double the amount of distance to complete the task compared to the long two-phase participants. We model this growth rate using a hierarchical Bayesian regression model (Detailed in Appendix I.3), with cumulative distance as the predictive variable. The experimental variables include interface type as a categorical variable, individual users modeled with random effects, and steps taken as

Difference and Effect Size of Distance per Option by Version



Figure 10: The figure shows the contrast distributions of the mean edit distance per option between pairwise experimental conditions, with the first row representing absolute differences and the second row depicting effect sizes. Main takeaway: is that participants in the long text estimated more edit distance per option compared to those in the short text and the long two-phase condition. Notably, the long two-phase interface required estimated only slightly more edit distances despite the longer survey length.



Figure 11: Edit Distance Per Action: This plot shows the frequency of specific edit distances at each step across the text interface and two-phase interface. Main takeaway: Participants in the long two-phase interface tend to make adjustments closer to their previous actions, resulting in visually less variance in edit distances throughout the entire survey.

a continuous variable. A truncated normal likelihood constrains cumulative distances to positive values and varies these distances across steps for each participant while masking incomplete data.

Figure 14 shows that the slope for the long text interface is approximately 4.7, meaning each step by the text interface would add 4.7 edit distance (94% HDI = [4.2, 5.4]), compared to the long two-phase interface, which shows a statistically significant difference with a mode of 1.4 (94% HDI = [1.3, 1.7]). These results explain that the variance in edit distance per action and the increase in per option edit distance are consistent across participants between the two groups, showing that the organization phase allows participants to focus on adjusting options within proximity without

having to navigate the interface to locate and make adjustments throughout the voting phase.

**Evidence from qualitative analysis:** Recall the differences in sources of cognitive load between the two experimental conditions: while two-phase interface participants make localized adjustments with nearby options, they experience cognitive demand from preference construction due to broader considerations that involve more options and higher-order values. Similarly, the qualitative results highlight that long text interface participants construct narrower preferences, yet their edit distance indicates broader movements across options.

Figure 12: Posterior variance differences (left) and effect sizes (right) in mean edit distance per step between text and two-phase interfaces for different survey lengths. Main takeaway: The long text interface had greater variance in edit distance per step, while differences in the short text condition were not statistically significant.



Figure 13: Cumulative edit distances over the survey for long text and long two-phase groups. Main takeaway: The long two-phase interface encourages smaller, incremental adjustments, leading to a flatter slope than the text interface.



Figure 14: Posterior distribution of slope differences (left) and effect sizes (right) in cumulative edit distance between interactive and two-phase interfaces for long QSs. Main takeaway: Participants in the interactive interface made larger adjustments compared to the two-phase interface.

Fewer long two-phase interface participants (60%, N=6) reported precise resource allocation as a source of demand compared to 90% in the text interface (N=9). We interpret this as former participants construct preliminary preferences during the organization phase, easing them to concentrate vote decisions as they focus more on deliberate preference building rather than mere completion. Conveniently positioning options with similar preferences reduced the need to look for an option and traverse the interface, allowing participants remain engaged in vote adjustments.

## 7 Clickstream data: Time participants spent



Figure 15: Total Time per Option. Each dot represents the time a participant took to complete an option, with the plot's shape showing the distribution within each group. The wider it is, the more dots there are. The three horizontal lines indicate the 25th, 50th, and 75th percentile annotated with value. The two-phase interface skewed slightly higher than the text interface Main takeaway: Two-phase interface participants spend longer time per option compared to its counterparts.
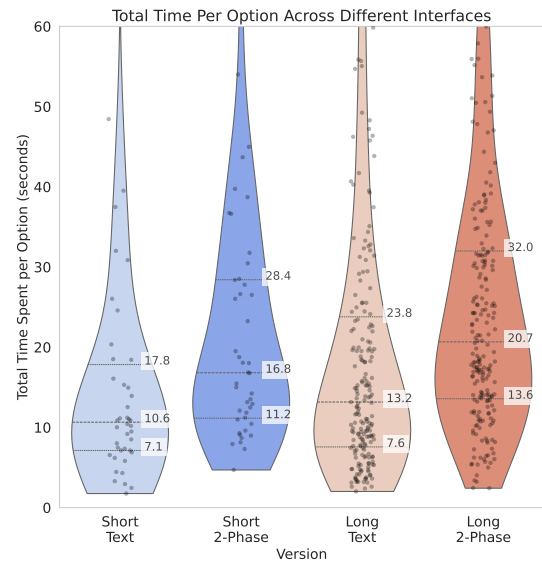
In addition to distance, participants in the short survey took an average of 2.7 minutes (short-text: $\mu$=2.3, $\sigma$=1.27; short two-phase: $\mu$=3, $\sigma$=1.02), while those in the long survey took 9.7 minutes (long-text: $\mu$=7.5, $\sigma$=3.45; long two-phase: $\mu$=11.95, $\sigma$=2.73). For a fairer comparison of interaction patterns, we analyze total **time-spend-per-option** using the QS system logs in this section. For participants in the two-phase interface conditions, this includes both organization and voting times for that option. The results are visualized in Figure 15.

Overall, participants spent slightly more time per option in the two-phase interface than in the text interface. To quantify these observations, we model the time data as predictive variables of separate Gamma distributions to characterize the continuous response times observed under distinct experimental conditions defined by survey length and interface type (Detailed in Appendix H). Each of the four resulting subsets of data is modeled independently, with

Figure 16: The figure shows the contrast distributions of the mean time to complete per option between pairwise experimental conditions, with the first row representing absolute differences and the second row depicting effect sizes. Main takeaway: is that participants in the long two-phase condition spent more time per option compared to those in the long text and short two-phase conditions. Additionally, short two-phase participants took longer per option than short text participants.

separate Gamma-distributed parameters governing the shape and rate of each group's time distributions.

We calculated the posterior differences between all pairwise comparisons of the four groups. The results in Figure 16 indicate that participants using the two-phase interface consistently spend more time per option than those using the text interface, regardless of survey length. For both the short and long QSs, participants most likely spend 6.1 seconds (94% HDI = [1.0, 11.0]) and 6.7 seconds (94% HDI = [3.7, 9.4]) more per option, respectively, with medium effect sizes of $d$=0.49 (94% HDI = [0.077, 0.89]) and $d$=0.41 (94% HDI = [0.24, 0.59]). In both cases, the intervals lie outside the ROPE of 0 ± 1, indicating statistical significance. These findings suggest that the two-phase interface encourages longer deliberation, particularly for long option surveys.

Some literature points out that increased time can lead to cognitive fatigue [38, 45], which can impair decision-making. Other decision science literature suggests that longer decision times can indicate deeper cognitive processing [15, 64]. Our qualitative analysis points to the latter.

Descriptively, participants in the long two-phase condition remained actively engaged during the voting phase, editing their votes an average of 39.3 times per participant ($\sigma$=39.3, range=19 − 63) compared to 39.1 times ($\sigma$=13.29, range=15 − 58) in the long text condition. This suggests that the two-phase interface does not reduce engagement despite the additional organization step.

Quantitatively, other than the difference in operational thinking and strategic consideration discussed in Section 5.3, we find that 37.5% of participants (N=15) who attribute time to *Decision Making* as a source of temporal demand frame such demand differently. We label a participant as *affirmative* if they describe the pressure to make decisions as a source of temporal demand. For example, S022 ✍ *So it didn't take too much time, but obviously there were a lot of things to consider, so there was some temporal demand.* is an affirmative statement. Conversely, we label a participant as *negative* if they express concern about the time and effort they have already

invested. For example, S024 ✍ *maybe I should just hurry up and make a decision.* is a negative statement.

50% of participants (N=5) in the long two-phase group describe the pressure to make decisions affirmatively and none negatively. This suggests that their pressure stems from having too many remaining decisions to make, rather than from the time already invested. This is reflected in their higher average time spent per option and overall time spent ($\mu$=716.86 seconds, $\sigma$=164.04 seconds) completing the QS compared to the long text group ($\mu$=449.64 seconds, $\sigma$=206.97 seconds). We interpret these results that participants are thoughtfully engaged in constructing their preferences and choose to invest additional time, rather than being driven by decision-related pressures or experiencing urgency.

Conversely, in the short text group, 50% of participants (N=5) express concern about the time and effort they have already invested ( S024 ✍ *maybe I should just hurry up and make a decision.* ) and none frame it affirmatively. Descriptively, participants in the short text group spend comparatively less time than those in the long QS (short text: $\mu$=139.83 seconds, $\sigma$=76.43 seconds; short two-phase: $\mu$=178.78 seconds, $\sigma$=61.07 seconds). This suggests that participants in the short text group expect themselves to complete the task sooner than they actually do.

Surprisingly, participants in the long text interface exhibit lower temporal demand compared to both the short text and long two-phase interfaces (Figure 17). Bayesian analysis (Appendix G.2.3) supports this finding, with posterior probabilities of 86.1% and 86.7%, respectively. This result is notable considering participants spent more time per option compared to those in the short text interface and traversed the longest distance among all three groups (Section 6). In addition, only 30% of participants (N=3) mention the time spent making a decision as a source of temporal demand. One possible explanation is that some participants are satisficing, as we pointed out in Section 5.4.

In summary, we interpret the result that participants in the two-phase interface spend more time per option as a sign of deeper

cognitive processing. This is further supported by examining participants' nuanced voting behaviors under budget constraint conditions for the long QS, which we omit here for brevity. Notably, two-phase interface participants make more small vote adjustments (i.e., adding or removing at most 2 votes on an option) when they have fewer remaining credits, further supporting our claim that they experience deeper engagement with preference construction, which we elaborate on further in Appendix F.



Figure 17: Temporal Demand Raw Score: Each dot represents a participant's subscale response. Main takeaway: Long text interface participants seem to express less temporal demand compared to the other experiment conditions.

## 8 Discussion and Future Work

In this section, we interpret our findings on cognitive load and respondent behavior in a QS. We highlight the rationale and elements behind the two-phase interface for preference construction and its potential to mitigate satisficing behaviors. We also offer usage and design recommendations for practitioners and outline future directions for improving QS interfaces.

### 8.1 Two-phase interface: a worthwhile trade-off

Survey designers seek thoughtful responses from participants. This means the interface should balance survey usability, respondent satisfaction, and the effort individuals invest in their responses. Our results indicate that the two-phase interface encouraged deeper participant engagement with the options and reduced satisficing behaviors, despite its increased time per option and higher cognitive load for the long QSs.

*8.1.1 Analysis through the lens of cognitive load theory.* Cognitive load theory [86], when applied to QSs, identifies three components of cognitive load: intrinsic load (the cognitive demand required to understand questions and response options), germane load (associated with deeper processing and preference evaluation), and extraneous load (stemming from navigating and operating the survey interface).

Participants were randomly assigned to experimental conditions, with survey lengths containing options randomly drawn from a common pool to control for intrinsic load within the same group.

When a QS is short, participants can engage with all options simultaneously. Participants using the two-phase interface traded a slightly longer survey response time for a potential reduction in cognitive load and edit distance. We interpret this as participants freeing up cognitive demand from extraneous load for germane load, prompting them to better construct and express their preferences.

When a QS is long, participants face more options, resulting in a higher intrinsic load at the start of the survey. We believe the two-phase interface traded longer survey response time and a potential increase in cognitive load for deeper engagement with the survey. This heightened cognitive load likely stemmed from making comparisons in both the organization and voting phases. Quantitatively, participants spent more time per option, suggesting deeper engagement while exerting limited extraneous load, as evidenced by shorter traversals during voting. Qualitatively, participants reported experiencing demand primarily from strategic considerations (germane load) rather than operational actions (extraneous load), which were common among text interface participants.

While some might argue that the additional organizing phase offers participants more opportunities to familiarize themselves with the options compared to text interface participants, the greater overall edit distance and high variance in edit distance per option suggest that text interface participants traversed the list frequently. This finding is further supported by qualitative data, where 70% of long-text participants (N=7) reported scanning the list while voting. This behavior suggests that while long-text participants had opportunities to familiarize themselves with the options, the explicit organization phase encouraged deeper reflection on their preferences.

The effect of the two-phase interface shows nuanced differences influencing cognitive load outcomes; however, both analyses suggest that the two-phase interface *shifted* participants' cognitive focus when completing QS.

*8.1.2 Potential in limiting Satisficing.* Qualitative findings (Section 5.4) on potential satisficing behavior highlight the importance of careful consideration when deploying a long QS. However, the two-phase interface appeared to limit satisficing behaviors, as evidenced by fewer observations compared to the long text interface for the long QS and none for the short QS. We believe the potential reasons lie in the design of the two-phase interface, which scaffolds the preference construction process.

The deliberate one-option-at-a-time presentation during the voting task in the two-phase interface reduced reliance on defaults and encouraged deeper reflection using cognitive strategies such as *problem decomposition* [81] and *dimension reduction*, both of which are known to reduce cognitive overload.

When asked about their experience with the interface, four participants highlighted how the organization phase supported their preference construction. S013 illustrated how the one-option-at-a-time approach reduced the dimensions of decision-making:

*[. . . ] it (organization phase) gives you time to just focus on that single thing and rank it based on how you feel at that moment.*

💬 S013 (S2P)

This focused mode enabled deeper reflection. When considering relative preferences among QS options, S009 described how it structurally decomposed the problem:

*[…] to have a preliminary categorization of all the topics […] (allowed me) to think about and process […] digest all the information prior to actually allocating the budget […]*

S009 (L2P)

This quote highlighted how participants' deliberation occurred during the organization phase, enabling them to focus on constructing preferences without worrying about budget management—both of which are cited sources of cognitive load. Although direct measurement of satisficing behavior reduction is challenging, qualitative data and participant feedback suggest that the two-phase interface potentially limits such behaviors. Based on this evidence, we recommend that long QSs be implemented with a two-phase interface and sufficient time for participants to complete the process. We suggest future research investigate the mental processes underlying satisficing behaviors in long QSs.

**In summary,** we argue that the trade-off of a longer completion time and potentially higher cognitive load in the two-phase interface is justified. Drawing on cognitive load theory, the interface fosters deeper engagement with the options. Additionally, our qualitative findings and participant feedback suggest that the interface may reduce satisficing, aligning with decision-makers' goals of obtaining thoughtful and deliberate responses from participants.

## 8.2 Preference Construction guided by Organize, Then Vote

Completing a QS involves a series of in-situ, difficult decision tasks as participants construct their preference over unfamiliar options [49], as one participant reflected:

*Oh, there are other aspects that I never care about. […] Why (should) I spend money on that?*

S037 (L2P)

We believe the two-phase interface supported participants' preference construction process when faced with unfamiliar options.

First, 40% of long-text participants (N=4) found it challenging to facilitate differentiation without organization tools that would allow grouping or drag-and-drop, as S025 said:

*I would like to be able to like, click and drag the categories themselves so I could maybe reorder them to like my priorities. […] make myself categories and subcategories out of this list …If I could organize it.*

S025 (LT)

In contrast, 60% (N=6) of long two-phase participants appreciated the upfront introduction of all options, which enabled them to organize and use drag-and-drop features to facilitate QS completion. Not only did participants use drag-and-drop options post-voting to reflect and ensure correct vote allocation, but drag-and-drop also enabled participants, like S039, to make fine-grained comparisons between options:

*I think the system was actually really helpful because I could just drag them. […] I can really compare them, I can drag this one up here, and then compare it to the top one […]*

S039 (S2P)

This supports our intention of applying Svenson [85]'s differentiation and consolidation theory, in which participants attempt to identify differences and eliminate less favorable options. The

organization phase and the drag-and-drop supported some degree of differentiation process.

*[…] the hardest part deciding in which category of place (prefernce bin) each issue is.*

S021 (L2P)

This quote by S021 best represents the potential of the organization phase in separating part of the difficult decisions one needs to make when differentiating their preferences during preference construction. With the selected options, the shorter edit distance of long two-phase interface participants suggested that they were consolidating their identified preferences through votes.

## 8.3 What We Learned: Quadratic Survey Usage and Design Recommendations

This study represents a crucial step toward developing better interfaces to support individuals responding to QSs by providing a deeper understanding of how survey respondents interact with QSs and the sources of cognitive load. In this subsection, we outline usage and design recommendations applicable to all applications of the quadratic mechanism.

*8.3.1 QS: Prioritizing Fewer Options or High-Stakes Evaluations.* We recommend deploying a QS with smaller sets of options or for critical evaluations, such as eliciting stakeholders' preferences before making investment decisions in hospital infrastructure. Our findings indicate that cognitive challenges and time requirements increase significantly as the number of options grows. For a long QS, while the two-phase interface helps mitigate some challenges, it does not eliminate them entirely, making adequate deliberation time essential. If a two-phase interface is unavailable, survey designers should present options in advance to allow participants to familiarize themselves and reflect before completing the QS.

*8.3.2 Facilitate Quadratic Mechanism Applications through Categorization, Not Ranking.* In a QS, the final ranking of preferences is typically a byproduct of vote allocation rather than a deliberate ranking effort. Participants did not explicitly rank options; instead, their preferences emerged dynamically through the voting process. To better support this preference construction, future quadratic mechanism interface designs should focus on helping participants categorize options effectively rather than ranking them directly. Facilitating differentiation among options is more critical than enabling precise manipulation for fine-tuning. We believe this approach extends beyond QSs to other ranking-based survey tools, such as ranked-choice voting and constant-sum surveys. Further research should examine how implementing such functionality influences survey respondents' mental models.

## 8.4 Future work: Opportunities for Better Budget Management

Budget management emerged as one of the participants' most prominent challenges, which the two-phase interface did not address. 35% of participants (N=14) emphasized that current quadratic mechanism applications support automated calculations, but noted their insufficiency. We identified three challenges for future work:

First, participants struggled to decide on an initial vote allocation. Some distributed credits equally across options, while others used

1, 2, or 3 votes as starting points. A few anchored their decisions to the tutorial's example of four upvotes. This suggests a need to better understand whether individuals have absolute value preferences among options. Second, 12.5% of participants (N=5) expressed confusion about the relationship between budget, votes, and outcomes, despite understanding their definitions. They struggled to make trade-offs between votes and budget, leading to frustration and hampered decision-making. Third, determining the absolute amount of credits in a QS is highly demanding. Designing interfaces and interactions to address the cold start challenge and help participants decide on the absolute vote value, while also considering ways to limit direct influences, remains an open question.

We believe that, with a well-designed interface backed by real-time computing and a better understanding of how individuals calculate trade-offs, we can provide innovative solutions to help participants more easily express their preferences using QSs.

## 9   Limitations

Evaluating the QS interface is challenging because of its novelty. We identified several limitations that warrant further research.

*Individual differences in cognitive capacity.* Variations in individual cognitive capacity influenced participants' performance and cognitive scores. For example, participants with greater experience in decision-making may be better able to manage multiple options. A within-subject study could clarify shifts in cognitive load, but deconstructing established preferences and altering options introduces additional complexity. Therefore, we opted for this in-depth, between-subject study, although the small sample size may introduce noise, potentially distorting the measurement of cognitive load. Future research should aim to quantify the impact of different QS interfaces on cognitive load at a larger scale. Furthermore, participants completed this study in a controlled laboratory environment, with options displayed on a large screen. Future work should also investigate how individuals respond to QSs on smaller devices and in less controlled environments.

*Limited experience with QSs.* Participants lacked prior experience with the QS interface. After completing a tutorial and quiz, participants proceeded to perform tasks using the QS interface. While participants understood the mechanics of QSs, their familiarity with the interface likely influenced their strategies and cognitive load. As quadratic mechanisms become more prevalent, future research could compare the performance of novices and experts.

*Limitations of Time and Distance as Proxies for Decision-Making Effort.* While time and distance are common metrics for quantifying the effort involved in decision-making, they do not capture without noise. Participants may have considered multiple options simultaneously. We acknowledge that these metrics are approximate indicators of decision-making effort. Despite these limitations, this approach provides valuable insights into decision-making within our experimental constraints.

*Other Limitations.* Finally, although we observe meaningful trends in the Bayesian statistical results, the small sample size limits our ability to establish statistical significance in cognitive load differences. Additionally, despite our best efforts to ensure transparency in the qualitative analysis, potential biases may have been introduced by relying on a single coder. Future work should address these limitations by incorporating larger sample sizes and multiple coders to enhance the reliability and generalizability of findings related to cognitive load in QSs.

## 10   Conclusion

This study introduces and evaluates a two-phase "Organize-then-Vote" interface to help QS respondents construct their preferences. We examined how the interface affected cognitive load and response behaviors across societal issues of varying lengths through an in-lab study, NASA-TLX, and interviews. The interface's organization and voting phases, designed to reduce cognitive overload by structuring the decision-making process, allowed respondents to differentiate between options before voting. Results revealed that the two-phase design reduced participants' edit distance between vote adjustments throughout the survey and they spent more time per option. Qualitative insights highlighted that the two-phase interface encouraged more iterative and reflective preference construction and its potential for reducing satisficing behaviors even though it did not clearly reduce the overall cognitive load for the longer QS. Nonetheless, this design shift promoted deeper engagement and strategic thinking compared to the text-based interface, by distributing cognitive effort more effectively. By integrating the organization and drag-and-drop functions, the interface facilitated both preference differentiation and consolidation, making it easier for respondents to refine their decisions. This two-phase interface design supports the development of future software tools that facilitate preference construction and promote the broader adoption of QSs. Future research should explore how to better support individuals' budget allocation and design interfaces for smaller devices.

## Acknowledgments

## References

[1] Duane F Alwin and Jon A Krosnick. 1985. The Measurement of Values in Surveys: A Comparison of Ratings and Rankings. *Public Opinion Quarterly* 49, 4 (1985), 535–552.

[2] Vitalik Buterin, Zoë Hitzig, and E. Glen Weyl. 2019. A Flexible Design for Funding Public Goods. *Management Science* 65, 11 (Nov. 2019), 5171–5187. doi:10.1287/mnsc.2019.3337

[3] Brad Cain. 2007. *A Review of the Mental Workload Literature.* Technical Report. Defense Technical Information Center.

[4] Charlotte Cavaillé, Daniel L. Chen, and Karine Van der Straeten. 2024. Who Cares? Measuring Differences in Preference Intensity. *Political Science Research and Methods* (2024), 1–17. doi:10.1017/psrm.2024.27

[5] Center for Civic Design. n.d.. Center for Civic Design. https://civicdesign.org/.

[6] Charity Navigator. 2023. Charity Navigator. https://www.charitynavigator.org/index.cfm?bay=search.categories.

[7] Ti-Chung Cheng, Tiffany Li, Yi-Hung Chou, Karrie Karahalios, and Hari Sundaram. 2021. "I Can Show What I Really like.": Eliciting Preferences via Quadratic Voting. *Proceedings of the ACM on Human-Computer Interaction* 5 (April 2021), 1–43. doi:10.1145/3449281
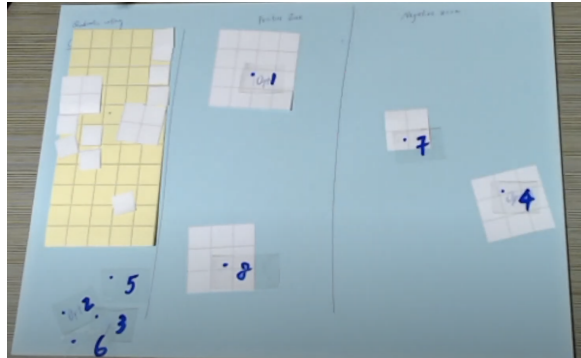
[8] Alexander Chernev, Ulf Böckenholt, and Joseph Goodman. 2015. Choice Overload: A Conceptual Review and Meta-Analysis. *Journal of Consumer Psychology* 25, 2 (April 2015), 333–358. doi:10.1016/j.jcps.2014.08.002

[9] Dana Chisnell. 2016. Democracy Is a Design Problem. *Journal of Usability Studies* 11, 3 (2016), 124–130.

[10] CivicDesign. 2015. Designing Usable Ballots Center for Civic Design. https://civicdesign.org/fieldguides/designing-usable-ballots/.

[11] Frederick G. Conrad, Benjamin B. Bederson, Brian Lewis, Emilia Peytcheva, Michael W. Traugott, Michael J. Hanmer, Paul S. Herrnson, and Richard G. Niemi. 2009. Electronic Voting Eliminates Hanging Chads but Introduces New Usability Challenges. *International Journal of Human-Computer Studies* 67, 1 (Jan. 2009), 111–124. doi:10.1016/j.ijhcs.2008.09.010

[12] Graham Cooper. 1998. Research into Cognitive Load Theory and Instructional Design at UNSW.

[13] M. P. Couper. 2001. Web Survey Design and Administration. *Public Opinion Quarterly* 65, 2 (2001), 230–253. doi:10.1086/322199

[14] Peter Coy. 2019. A New Way of Voting That Makes Zealotry Expensive - Bloomberg. *Bloomberg* (May 2019).

[15] Kahneman Daniel. 2017. *Thinking, Fast and Slow.* Farrar, Straus and Giroux.

[16] Shaneé Dawkins, Tony Sullivan, Greg Rogers, E. Vincent Cross, Lauren Hamilton, and Juan E. Gilbert. 2009. Prime III: An Innovative Electronic Voting Interface. In *Proceedings of the 14th International Conference on Intelligent User Interfaces (IUI '09).* Association for Computing Machinery, New York, NY, USA, 485–486. doi:10.1145/1502650.1502727

[17] Antonio Drommi, Gregory W Ulferts, and Dan Shoemaker. 2001. Interface Design: A Focus on Cognitive Science. In *The Proceedings of ISECON 2001,* Vol. 18.

[18] E. Glen Weyl, Kasia Sitkiewicz, and Petar Maymounkov. 2023. Gov4git: A Decentralized Platform for Community Governance.

[19] Erik J Engstrom and Jason M Roberts. 2020. *The Politics of Ballot Design: How States Shape American Democracy.* Cambridge University Press.

[20] Sarah P. Everett, Kristen K. Greene, Michael D. Byrne, Dan S. Wallach, Kyle Derr, Daniel Sandler, and Ted Torous. 2008. Electronic Voting Machines versus Traditional Methods: Improved Preference, Similar Performance. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems (CHI '08).* Association for Computing Machinery, New York, NY, USA, 883–892. doi:10.1145/1357054.1357195

[21] Habiba Farzand, David Al Baiaty Suarez, Thomas Goodge, Shaun Alexander Macdonald, Karola Marky, Mohamed Khamis, and Paul Cairns. 2024. Beyond Aesthetics: Evaluating Response Widgets for Reliability & Construct Validity of Scale Questionnaires. In *Extended Abstracts of the 2024 CHI Conference on Human Factors in Computing Systems (CHI EA '24).* Association for Computing Machinery, New York, NY, USA, 1–7. doi:10.1145/3613905.3650751

[22] N. T. Feather. 1973. The Measurement of Values: Effects of Different Assessment Procedures. *Australian Journal of Psychology* 25, 3 (Dec. 1973), 221–231. doi:10.1080/00049537308255849

[23] Martin Feick, Niko Kleer, Anthony Tang, and Antonio Krüger. 2020. The Virtual Reality Questionnaire Toolkit. In *Adjunct Proceedings of the 33rd Annual ACM Symposium on User Interface Software and Technology.* 68–69.

[24] Robert Ferber. 1952. Order Bias in a Mail Survey. *Journal of Marketing* 17, 2 (1952), 171–178. doi:10.2307/1248043 jstor:1248043

[25] Eibe Frank and Mark Hall. 2001. A Simple Approach to Ordinal Classification. In *Machine Learning: ECML 2001: 12th European Conference on Machine Learning Freiburg, Germany, September 5–7, 2001 Proceedings 12.* Springer, 145–156.

[26] Luis Mota Freitas and Wilfredo L. Maldonado. 2024. Quadratic Funding with Incomplete Information. *Social Choice and Welfare* (Feb. 2024). doi:10.1007/s00355-024-01512-7

[27] Qin Gao, Yang Wang, Fei Song, Zhizhong Li, and Xiaolu Dong. 2013. Mental Workload Measurement for Emergency Operating Procedures in Digital Nuclear Power Plants. *Ergonomics* 56, 7 (July 2013), 1070–1085. doi:10.1080/00140139.2013.790483

[28] Gerd Gigerenzer and Daniel G. Goldstein. 1996. Reasoning the Fast and Frugal Way: Models of Bounded Rationality. *Psychological Review* 103, 4 (1996), 650–669. doi:10.1037/0033-295X.103.4.650

[29] Juan E. Gilbert, Jerone Dunbar, Alvitta Ottley, and John Mark Smotherman. 2013. Anomaly Detection in Electronic Voting Systems. *Information Design Journal (IDJ)* 20, 3 (Sept. 2013), 194–206. doi:10.1075/idj.20.3.01gil

[30] Gitcoin. [n. d.]. Read the Whitepaper | Gitcoin. https://www.gitcoin.co/whitepaper/read.

[31] Theodore Groves and John Ledyard. 1977. Optimal Allocation of Public Goods: A Solution to the "Free Rider" Problem. *Econometrica* 45, 4 (1977), 783–809. doi:10.2307/1912672 jstor:1912672

[32] Eija Haapalainen, SeungJun Kim, Jodi F. Forlizzi, and Anind K. Dey. 2010. Psycho-Physiological Measures for Assessing Cognitive Load. In *Proceedings of the 12th ACM International Conference on Ubiquitous Computing.* ACM, Copenhagen Denmark, 301–310. doi:10.1145/1864349.1864395

[33] Sandra G. Hart. 2006. Nasa-Task Load Index (NASA-TLX); 20 Years Later. *Proceedings of the Human Factors and Ergonomics Society Annual Meeting* 50, 9

[34] Sandra G Hart and Lowell E Staveland. 1988. Development of NASA-TLX (Task Load Index): Results of Empirical and Theoretical Research. In *Advances in Psychology.* Vol. 52. Elsevier, 139–183.

[35] Susan C. Herring and Ashley R. Dainas. 2020. Gender and Age Influences on Interpretation of Emoji Functions. *ACM Transactions on Social Computing* 3, 2 (June 2020), 1–26. doi:10.1145/3375629

[36] Emin İbili. 2019. Effect of Augmented Reality Environments on Cognitive Load: Pedagogical Effect, Instructional Design, Motivation and Interaction Interfaces. *International Journal of Progressive Education* 15, 5 (2019), 42–57.

[37] Sheena S Iyengar and Mark R Lepper. 2000. When Choice Is Demotivating: Can One Desire Too Much of a Good Thing? *Journal of personality and social psychology* 79, 6 (2000), 995.

[38] Enamul Karim, Hamza Reza Pavel, Sama Nikanfar, Aref Hebri, Ayon Roy, Harish Ram Nambiappan, Ashish Jaiswal, Glenn R Wylie, and Fillia Makedon. 2024. Examining the Landscape of Cognitive Fatigue Detection: A Comprehensive Survey. *Technologies* 12, 3 (2024), 38.

[39] Matthew Kay, Gregory L Nelson, and Eric B Hekler. 2016. Researcher-Centered Design of Statistics: Why Bayesian Statistics Better Fit the Culture and Incentives of HCI. In *Proceedings of the 2016 CHI Conference on Human Factors in Computing Systems.* 4521–4532.

[40] Aman Khullar, Priyadarshi Hitesh, Shoaib Rahman, Deepak Kumar, Rachit Pandey, Praveen Kumar, Rajeshwari Tripathi, Prince Prince, Ankit Akash Jha, Himanshu Himanshu, and Aaditeshwar Seth. 2021. Costs and Benefits of Conducting Voice-Based Surveys versus Keypress-Based Surveys on Interactive Voice Response Systems. In *Proceedings of the 4th ACM SIGCAS Conference on Computing and Sustainable Societies (Compass '21).* Association for Computing Machinery, New York, NY, USA, 288–298. doi:10.1145/3460112.3471963

[41] N. D. Kieruj and G. Moors. 2010. Variations in Response Style Behavior by Response Scale Format in Attitude Research. *International Journal of Public Opinion Research* 22, 3 (Sept. 2010), 320–342. doi:10.1093/ijpor/edq001

[42] Soomin Kim, Jinsu Eun, Joseph Seering, and Joonhwan Lee. 2021. Moderator Chatbot for Deliberative Discussion: Effects of Discussion Structure and Discussant Facilitation. *Proc. ACM Hum.-Comput. Interact.* 5, CSCW1, Article 87 (April 2021). doi:10.1145/3449161

[43] Soomin Kim, Joonhwan Lee, and Gahgene Gweon. 2019. Comparing Data from Chatbot and Web Surveys: Effects of Platform and Conversational Style on Survey Response Quality. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems.* ACM, Glasgow Scotland Uk, 1–12. doi:10.1145/3290605.3300316

[44] Jon A Krosnick, Charles M Judd, and Bernd Wittenbrink. 2018. The Measurement of Attitudes. In *The Handbook of Attitudes.* Routledge, 45–105.

[45] Thomas Kundinger, Celena Mayr, and Andreas Riener. 2020. Towards a Reliable Ground Truth for Drowsiness: A Complexity Analysis on the Example of Driver Fatigue. *Proceedings of the ACM on Human-Computer Interaction* 4, EICS (June 2020), 1–18. doi:10.1145/3394980

[46] Benjamin Lafreniere, Andrea Bunt, and Michael Terry. 2014. Task-Centric Interfaces for Feature-Rich Software. In *Proceedings of the 26th Australian Computer-Human Interaction Conference on Designing Futures: The Future of Design (OzCHI '14).* Association for Computing Machinery, New York, NY, USA, 49–58. doi:10.1145/2686612.2686620

[47] Steven P Lalley, E Glen Weyl, et al. 2016. Quadratic Voting. *Available at SSRN* (2016).

[48] Seunghyun "Tina" Lee, Yilin Elaine Liu, Ljilja Ruzic, and Jon Sanford. 2016. Universal Design Ballot Interfaces on Voting Performance and Satisfaction of Voters with and without Vision Loss. In *Proceedings of the 2016 CHI Conference on Human Factors in Computing Systems (CHI '16).* Association for Computing Machinery, New York, NY, USA, 4861–4871. doi:10.1145/2858036.2858567

[49] Sarah Lichtenstein and Paul Slovic (Eds.). 2006. *The Construction of Preference* (1. publ ed.). Cambridge University Press, Cambridge.

[50] Michael Xieyang Liu, Aniket Kittur, and Brad A. Myers. 2021. To Reuse or Not to Reuse? A Framework and System for Evaluating Summarized Knowledge. *Proc. ACM Hum.-Comput. Interact.* 5, CSCW1, Article 166 (April 2021). doi:10.1145/3449240

[51] Richard McElreath. 2018. *Statistical Rethinking: A Bayesian Course with Examples in R and Stan.* Chapman and Hall/CRC.

[52] George A. Miller. 1956. The Magical Number Seven, plus or Minus Two: Some Limits on Our Capacity for Processing Information. *Psychological Review* 63, 2 (1956), 81–97. doi:10.1037/h0043158

[53] Henry Montgomery. 1983. Decision Rules and the Search for a Dominance Structure: Towards a Process Model of Decision Making. In *Advances in Psychology.* Vol. 14. Elsevier, 343–369. doi:10.1016/S0166-4115(08)62243-8

[54] William F. Moroney and Joyce A. Cameron. 2019. *Questionnaire Design: How to Ask the Right Questions of the Right People at the Right Time to Get the Information You Need.* Human Factors and Ergonomics Society.

[55] Robert Münscher, Max Vetter, and Thomas Scheuerle. 2016. A Review and Taxonomy of Choice Architecture Techniques. *Journal of Behavioral Decision*

(Oct. 2006), 904–908. doi:10.1177/154193120605000909

*Making* 29, 5 (2016), 511–524. doi:10.1002/bdm.1897

[56] Ryan Naylor et al. 2017. First Year Student Conceptions of Success: What Really Matters? *Student Success* 8, 2 (2017), 9–19.

[57] Syavash Nobarany, Louise Oram, Vasanth Kumar Rajendran, Chi-Hsiang Chen, Joanna McGrenere, and Tamara Munzner. 2012. The Design Space of Opinion Measurement Interfaces: Exploring Recall Support for Rating and Ranking. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems (Chi ’12)*. Association for Computing Machinery, New York, NY, USA, 2035–2044. doi:10.1145/2207676.2208351

[58] A Norman Donald. 2013. *The Design of Everyday Things.* MIT Press.

[59] Judith S. Olson and Wendy A. Kellogg (Eds.). 2014. *Ways of Knowing in HCI.* Springer, New York, NY. doi:10.1007/978-1-4939-0378-8

[60] Sharon Oviatt. 2006. Human-Centered Design Meets Cognitive Load Theory: Designing Interfaces That Help People Think. In *Proceedings of the 14th ACM International Conference on Multimedia (Mm ’06)*. Association for Computing Machinery, New York, NY, USA, 871–880. doi:10.1145/1180639.1180831

[61] Oskar Palinko, Andrew L. Kun, Alexander Shyrokov, and Peter Heeman. 2010. Estimating Cognitive Load Using Remote Eye Tracking in a Driving Simulator. In *Proceedings of the 2010 Symposium on Eye-Tracking Research & Applications - ETRA ’10*. ACM Press, Austin, Texas, 141. doi:10.1145/1743666.1743701

[62] Christian Jilek Paula Gauselmann, Yannick Runge and Tobias Tempel. 2023. A Relief from Mental Overload in a Digitalized World: How Context-Sensitive User Interfaces Can Enhance Cognitive Performance. *International Journal of Human–Computer Interaction* 39, 1 (2023), 140–150. doi:10.1080/10447318.2022.2041882 arXiv:https://doi.org/10.1080/10447318.2022.2041882

[63] John W. Payne, James R. Bettman, and Eric J. Johnson. 1988. Adaptive Strategy Selection in Decision Making. *Journal of Experimental Psychology: Learning, Memory, and Cognition* 14, 3 (July 1988), 534–552. doi:10.1037/0278-7393.14.3.534

[64] John W. Payne, James R. Bettman, and Eric J. Johnson. 1993. *The Adaptive Decision Maker.* Cambridge University Press, Cambridge. doi:10.1017/CBO9781139173933

[65] PDSI. 2022. Quadratic Voting Frontend. Public Digital Innovation Space.

[66] Martin Pielot and Mario Callegaro. 2024. Did You Misclick? Reversing 5-Point Satisfaction Scales Causes Unintended Responses. In *Proceedings of the CHI Conference on Human Factors in Computing Systems*. ACM, Honolulu HI USA, 1–7. doi:10.1145/3613904.3642397

[67] Eric A Posner and E Glen Weyl. 2018. *Radical Markets: Uprooting Capitalism and Democracy for a Just Society.* Princeton University Press.

[68] David Quarfoot, Douglas von Kohorn, Kevin Slavin, Rory Sutherland, David Goldstein, and Ellen Konar. 2017. Quadratic Voting in the Wild: Real People, Real Votes. *Public Choice* 172, 1-2 (2017), 283–303.

[69] Whitney Quesenbery. 2020. Opinion | Good Design Is the Secret to Better Democracy. *The New York Times* (Oct. 2020).

[70] RadicalxChange. [n. d.]. About. https://www.radicalxchange.org/wiki/about.

[71] RadicalxChange Foundation. 2024. Quadratic Voting.

[72] Helena M Reis, Simone S Borges, Vinicius HS Durelli, Luis Fernando de S Moro, Anarosa AF Brandao, Ellen F Barbosa, Leônidas O Brandao, Seiji Isotani, Patricia A Jaques, and Ig I Bittencourt. 2012. Towards Reducing Cognitive Load and Enhancing Usability through a Reduced Graphical User Interface for a Dynamic Geometry System: An Experimental Study. In *2012 IEEE International Symposium on Multimedia*. IEEE, 445–450.

[73] Duncan Rintoul. [n. d.]. Visual and Animated Response Formats in Web Surveys: Do They Produce Better Data, or Is It All Just Fun and Games? ([n. d.]), 126.

[74] Adam Rogers. 2019. Colorado Tried a New Way to Vote: Make People Pay—Quadratically | WIRED. *Wired* (April 2019).

[75] Susana Rubio, Eva Díaz, Jesús Martín, and José M. Puente. 2004. Evaluation of Subjective Mental Workload: A Comparison of SWAT, NASA-TLX, and Workload Profile Methods. *Applied Psychology* 53, 1 (2004), 61–86. doi:10.1111/j.1464-0597.2004.00161.x

[76] Thomas L. Saaty. 1987. Principles of the Analytic Hierarchy Process. In *Expert Judgment and Expert Systems*, Jeryl L. Mumpower, Ortwin Renn, Lawrence D. Phillips, and V. R. R. Uppuluri (Eds.). Springer, Berlin, Heidelberg, 27–73. doi:10.1007/978-3-642-86679-1_3

[77] Thomas L Saaty and Mujgan S Ozdemir. 2003. Why the Magic Number Seven plus or Minus Two. *Mathematical and computer modelling* 38, 3-4 (2003), 233–244.

[78] Stoo Sepp, Steven J. Howard, Sharon Tindall-Ford, Shirley Agostinho, and Fred Paas. 2019. Cognitive Load Theory and Human Movement: Towards an Integrated Model of Working Memory. *Educational Psychology Review* 31, 2 (June 2019), 293–317. doi:10.1007/s10648-019-09461-9

[79] Anuj K. Shah, Eldar Shafir, and Sendhil Mullainathan. 2015. Scarcity Frames Value. *Psychological Science* 26, 4 (2015), 402–412.

[80] Herbert A. Simon. 1955. A Behavioral Model of Rational Choice. *The Quarterly Journal of Economics* 69, 1 (1955), 99–118. doi:10.2307/1884852 jstor:1884852

[81] Herbert A. Simon. 1996. *The Sciences of the Artificial* (3rd ed ed.). MIT Press, Cambridge, Mass.

[82] Tobin South, Leon Erichsen, Shrey Jain, Petar Maymounkov, Scott Moore, and E. Glen Weyl. 2024. Plural Management. doi:10.2139/ssrn.4688040

[83] Fritz Strack and Leonard L. Martin. 1987. Thinking, Judging, and Communicating: A Process Account of Context Effects in Attitude Surveys. In *Social Information Processing and Survey Methodology*, Hans-J. Hippler, Norbert Schwarz, and Seymour Sudman (Eds.). Springer, New York, NY, 123–148. doi:10.1007/978-1-4612-4798-2_7

[84] Kathryn Summers, Dana Chisnell, Drew Davies, Noel Alton, and Megan McKeever. 2014. Making Voting Accessible: Designing Digital Ballot Marking for People with Low Literacy and Mild Cognitive Disabilities. In *2014 Electronic Voting Technology Workshop/Workshop on Trustworthy Elections (EVT/WOTE 14)*.

[85] Ola Svenson. 1992. Differentiation and Consolidation Theory of Human Decision Making: A Frame of Reference for the Study of Pre- and Post-Decision Processes. *Acta Psychologica* 80, 1-3 (Aug. 1992), 143–168. doi:10.1016/0001-6918(92)90044-E

[86] John Sweller. 2011. Cognitive Load Theory. In *Psychology of Learning and Motivation*. Vol. 55. Elsevier, 37–76. doi:10.1016/B978-0-12-387691-1.00002-8

[87] Taipei Representative Office in the U.K. 2024. Taiwan Digital Minister Highlights Country’s Use of Technology to Bolster Democracy in FT Interview.

[88] Richard H. Thaler and Cass R. Sunstein. 2008. *Nudge: Improving Decisions about Health, Wealth, and Happiness.* Yale University Press, New Haven, CT, US. x, 293 pages.

[89] Jerry P Timbrook. 2013. *A Comparison of a Traditional Ranking Format to a Drag-and-Drop Format with Stacking.* Ph. D. Dissertation. University of Dayton.

[90] Vera Toepoel and Frederik Funke. 2018. Sliders, Visual Analogue Scales, or Buttons: Influence of Formats and Scales in Mobile and Desktop Surveys. *Mathematical Population Studies* 25, 2 (April 2018), 112–122. doi:10.1080/08898480.2018.1439245

[91] Vera Toepoel, Brenda Vermeeren, and Baran Metin. 2019. Smileys, Stars, Hearts, Buttons, Tiles or Grids: Influence of Response Format on Substantive Response, Questionnaire Experience and Response Time. *Bulletin of Sociological Methodology/Bulletin de Méthodologie Sociologique* 142, 1 (April 2019), 57–74. doi:10.1177/0759106319834665

[92] Amos Tversky and Daniel Kahneman. 1982. Judgments of and by Representativeness. In *Judgment under Uncertainty: Heuristics and Biases*, Daniel Kahneman, Paul Slovic, and AmosEditors Tversky (Eds.). Cambridge University Press, Cambridge, 84–98.

[93] Muhsin Ugur, Dvijesh Shastri, Panagiotis Tsiamyrtzis, Malcolm Dcosta, Allison Kalpakci, Carla Sharp, and Ioannis Pavlidis. 2015. Evaluating Smartphone-Based User Interface Designs for a 2d Psychological Questionnaire. In *Proceedings of the 2015 ACM International Joint Conference on Pervasive and Ubiquitous Computing*. 275–282.

[94] Paul Van Schaik and Jonathan Ling. 2007. Design Parameters of Rating Scales for Web Sites. *ACM Transactions on Computer-Human Interaction (TOCHI)* 14, 1 (2007), 4–es.

[95] Jonathan N. Wand, Kenneth W. Shotts, Jasjeet S. Sekhon, Walter R. Mebane, Michael C. Herron, and Henry E. Brady. 2001. The Butterfly Did It: The Aberrant Vote for Buchanan in Palm Beach County, Florida. *The American Political Science Review* 95, 4 (2001), 793–810. jstor:3117714

[96] Jing Wei, Weiwei Jiang, Chaofan Wang, Difeng Yu, Jorge Goncalves, Tilman Dingler, and Vassilis Kostakos. 2022. Understanding How to Administer Voice Surveys through Smart Speakers. *Proc. ACM Hum.-Comput. Interact.* 6, CSCW2, Article 548 (Nov. 2022), 548 pages. doi:10.1145/3555606

[97] Bert Weijters, Elke Cabooter, and Niels Schillewaert. 2010. The Effect of Rating Scale Format on Response Styles: The Number of Response Categories and Response Category Labels. *International Journal of Research in Marketing* 27, 3 (Sept. 2010), 236–247. doi:10.1016/j.ijresmar.2010.02.004

[98] Bert Weijters, Kobe Millet, and Elke Cabooter. 2021. Extremity in Horizontal and Vertical Likert Scale Format Responses. Some Evidence on How Visual Distance between Response Categories Influences Extreme Responding. *International Journal of Research in Marketing* 38, 1 (March 2021), 85–103. doi:10.1016/j.ijresmar.2020.04.002

[99] Christopher D Wickens and Anthony D Andre. 1990. Proximity Compatibility and Information Display: Effects of Color, Space, and Objectness on Information Integration. *Human factors* 32, 1 (1990), 61–77.

[100] yehjxraymond. 2024. Yehjxraymond/Qv-App.

[101] Amy X. Zhang and Justin Cranshaw. 2018. Making Sense of Group Chat through Collaborative Tagging and Summarization. *Proc. ACM Hum.-Comput. Interact.* 2, CSCW, Article 196 (Nov. 2018), 196 pages. doi:10.1145/3274465
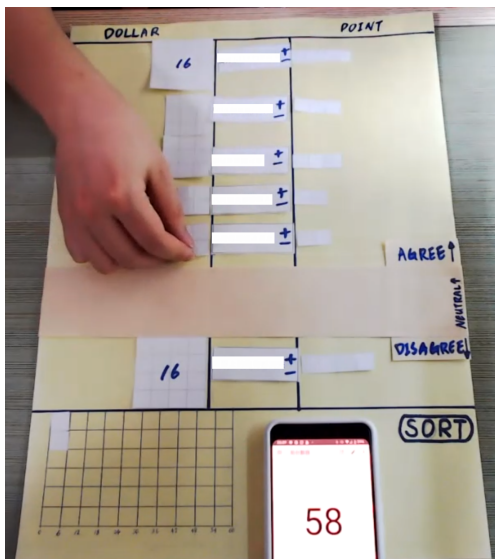
# A  Interface design process

In this section, we outline the design process leading to our final interface.



(a) In this paper prototype, issues are denoted by different numbers that appear on mouseover. Pretest respondents can move options anywhere in the two sections of the interface, one denoting positive and one negative. The blocks represent the cost for each option, with no indication of the number of current votes. The credits are shown in the yellow box on the left.



(b) This paper prototype separates the positive and negative areas with a 'band' at the center. Undecided options are placed inside this band. The cost and the votes on both sides of the interface are denoted by small blocks. The budget is shown in the yellow box below the interface with a numerical counter.

Figure 18: Initial paper prototypes designed for QS interface.

## A.1  Prototype 1: Ranking-Vote

Our first prototype emerged after various paper prototypes, such as those shown in Figure 18. Through pre-testing, we observed that participants engaging with QS needed interface support for organizing options and managing their credits. In this study, we decided to focus on the former.

Since participants needed to position options within the interface, and the end result formed a ranked list, we tested whether ranking options before voting would help establish an individual's relative preferences in Prototype 1 ( Figure 19). This prototype allowed respondents to reposition options before voting. However, pre-test respondents rarely moved the options and questioned the necessity of a full ranking, as it did not influence their QS submission. Additionally, many were unaware that the options were draggable. These findings suggest that a full ranking is unnecessary for establishing relative preferences. Therefore, we decided to ask respondents to select a subset of options rather than requiring a full ranking of all options.
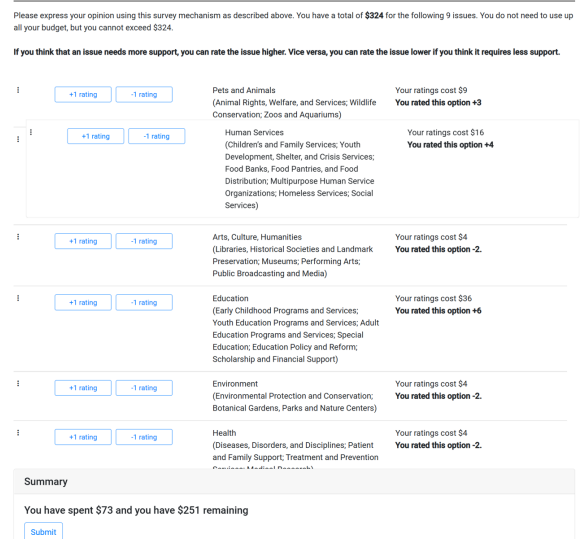


Figure 19: A Ranking-Vote Prototype: This prototype tests whether ranking options prior to voting helps establish an individual's relative preferences. Each option is draggable, allowing users to position it within the full list of options. Votes can be adjusted using the buttons on the left side of the interface, while the vote count and costs are displayed on the right. A summary box remains fixed at the bottom of the screen for easy reference.

## A.2  Prototype 2: Select-then-Vote

Based on feedback from Prototype 1, instead of *allowing* individuals to rank options, Prototype 2 implemented a two-phase process that *intentionally* asks respondents to select options to express opinions before voting.

As shown in Figure 20, survey respondents selected their preferred options (Figure 20a), and the interface positioned these options at the top of the list for voting (Figure 20b). We identified several issues during the prototype 2 pretest: many respondents marked most options as 'options they care about,' which undermined the design's purpose. Additionally, the lack of clear distinction between selected and unselected options confused respondents about the necessity of Step 1. Thus, we need a clearer distinction

**(a) Options are dragged and dropped to the 'Option You Care About' box.**



**(b) The previous step collapses showing all voting options.**

**Figure 20: A Select-then-Vote Prototype: The goal of this prototype is to nudge participants to focus on a subset of options to vote, rather than ranking all of them. This prototype introduces a two-step voting process. As shown in Fig. 20a, the first step involves selecting options for further consideration. Important options are placed at the top of the list for voting shown in Fig. 20b, but options can be placed anywhere on the list if desired. The rest of the controls follows the previous prototype.**



**(a) The Organization Interface: Options are shown initially in the first bin labeled as 'I don't know.' Survey respondents can then drag and drop these options into the latter bins: Lean Positive, Lean Neutral, or Lean Negative. Only the details of each option are shown on this interface.**



**(b) The Voting Interface: Voting controls appear on the left side of each option, showing the current votes and associated costs on the right. A budget summary sticks to the bottom of the screen.**

**Figure 21: Organize-then-Vote Prototype: The goal of this prototype is to encourage participants to derive finer grain categories among options before voting. Survey respondents first organize their thoughts into categories and then vote on the options.**

and connection between the two phases to effectively construct relative preferences.

## A.3 Prototype 3: Organize-then-Vote

Figure 21 shows the final prototype, which builds on our previous takeaway by introducing finer-grained groupings and establishing a clearer connection between option organization and voting position. Specifically, we provided three categories: Lean Positive, Lean Negative, and Lean Neutral. Initially, respondents see all options listed under a section labeled 'I don't know,' which displays only the option descriptions and not the vote controls. They are then asked to move these options into one of the three categories. On the subsequent page, voting controls and additional information appear for each option, reinforcing the connection between option grouping, position, and voting controls.

Feedback indicated that survey respondents are comfortable with the two-phase organize-then-vote design, demonstrating it as a central strategy for our interface development. However, we identified several areas for enhancement: First, the dragging and dropping mechanism in the organization phase is cumbersome and may inadvertently suggest a ranking process, contrary to our intentions. Second, placing unorganized options at the top of the voting list is counterintuitive. Third, the voting controls are disconnected from the option summaries, dividing attention between the left and right sides of the screen. These insights guided refinements in the final two-phase interface, adhering to the organize-then-vote framework.

## B Voting Interface Breakdown

In this section, we outline additional literature that informed this study. There are two sets of literature that we surveyed: Survey response format and voting interfaces.

## B.1 Survey response format

Research in the marketing and research communities focusing on survey and questionnaire design, usability, and interactions examines the influence of presentation styles and 'response format.' Weijters et al. [98] demonstrated that horizontal distances between options are more influential than vertical distances, with the latter recommended for reduced bias. Slider bars, which operate on a drag-and-drop principle, show lower mean scores and higher non-response rates compared to buttons, indicating they are more prone to bias and difficult to use. In contrast, visual analog scales that operate on a point-and-click principle perform better [90]. These studies show how even small design changes can have a large impact on usability, highlighting the importance of designing interfaces that prioritize human-centered interaction rather than focusing solely on functionality.

## B.2 Voting Interfaces

Compared to digital survey interfaces, voting interfaces are a specialized type of survey interface can significantly influence democratic processes [9, 10, 19] and often have consequential impacts. We categorize these related works into three main categories detailed below:

*Designs that shifted voter decisions:* For example, states without straight-party ticket voting (where voters can select all candidates from one party through a single choice) exhibited higher rates of split-ticket voting [19]. Another example from the Australian ballot showing incumbency advantages is where candidates are listed by the office they are running for, with no party labels or boxes.

*Designs that influenced errors:* Butterfly ballots, an atypical design, may have influenced the outcome of the 2000 U.S. Presidential Election [95]. It increased voter errors because voters could not correctly identify the punch hole on the ballot. Splitting contestants across columns increases the chance for voters to overvote [69]. On the other hand, Everett et al. [20] showed the use of incorporating physical voting behaviors, like lever voting, into graphical user interfaces increased satisfaction while maintaining efficiency and effectiveness.

*Designs that incorporated technologies:* Other projects like the Caltech-MIT Voting Technology Project addresses accessibility challenges, resulting in innovations like EZ Ballot [48], Anywhere Ballot [84], and Prime III [16]. In addition, Gilbert et al. [29] investigated optimal touchpoints on voting interfaces, and Conrad et al. [11] examined zoomable voting interfaces.

Response format literature and voting interfaces informed how interfaces significantly influence respondent behavior, decision accuracy, and cognitive load. These burdens are especially problematic for complex systems like QS, where high cognitive demands may deter researchers and users alike. Developing effective, human-centered interfaces for QS could enhance usability, reduce cognitive overload, and increase adoption in both research and practical applications.

## C List of Options

We provide the full list of options presented on the survey.

- **Animal Rights, Welfare, and Services:** Protect animals from cruelty, exploitation and other abuses, provide veterinary services and train guide dogs.
- **Wildlife Conservation:** Protect wildlife habitats, including fish, wildlife, and bird refuges and sanctuaries.
- **Zoos and Aquariums:** Support and invest in zoos, aquariums and zoological societies in communities throughout the country.
- **Libraries, Historical Societies and Landmark Preservation:** Support and invest public and specialized libraries, historical societies, historical preservation programs, and historical estates.
- **Museums:** Support and invest in maintaining collections and provide training to practitioners in traditional arts, science, technology, and natural history.
- **Performing Arts:** Support symphonies, orchestras, and other musical groups; ballets and operas; theater groups; arts festivals; and performance halls and cultural centers.
- **Public Broadcasting and Media:** Support public television and radio stations and networks, as well as providing other independent media and communications services to the public.
- **Community Foundations:** Promote giving by managing long-term donor-advised charitable funds for individual givers and distributing those funds to community-based charities over time.

- **Housing and Neighborhood Development:** Lead and finance development projects that invest in and improve communities by providing utility assistance, small business support programs, and other revitalization projects.
- **Jewish Federations:** Focus on a specific geographic region and primarily support Jewish-oriented programs, organizations and activities through grantmaking efforts
- **United Ways:** Identify and resolve community issues through partnerships with schools, government agencies, businesses, and others, with a focus on education, income and health.
- **Adult Education Programs and Services:** Provide opportunities for adults to expand their knowledge in a particular field or discipline, learn English as a second language, or complete their high school education.
- **Early Childhood Programs and Services:** Provide foundation-level learning and literacy for children prior to entering the formal school setting.
- **Education Policy and Reform:** Promote and provide research, policy, and reform of the management of educational institutions, educational systems, and education policy.
- **Scholarship and Financial Support:** Support and enable students to obtain the financial assistance they require to meet their educational and living expenses while in school.
- **Special Education:** Provide services, including placement, programming, instruction, and support for gifted children and youth or those with disabilities requiring modified curricula, teaching methods, or materials.
- **Youth Education Programs and Services:** Provide programming, classroom instruction, and support for school-aged students in various disciplines such as art education, STEM, outward bound learning experiences, and other programs that enhance formal education.
- **Botanical Gardens, Parks, and Nature Centers:** Promote preservation and appreciation of the environment, as well as leading anti-litter, tree planting and other environmental beautification campaigns.
- **Environmental Protection and Conservation:** Develop strategies to combat pollution, promote conservation and sustainable management of land, water, and energy resources, protect land, and improve the efficiency of energy and waste material usage.
- **Diseases, Disorders, and Disciplines:** Seek cures for diseases and disorders or promote specific medical disciplines by providing direct services, advocating for public support and understanding, and supporting targeted medical research.
- **Medical Research:** Devote and invest in efforts on researching causes and cures of disease and developing new treatments.
- **Patient and Family Support:** Support programs and services for family members and patients that are diagnosed with a serious illness, including wish granting programs, camping programs, housing or travel assistance.
- **Treatment and Prevention Services:** Provide direct medical services and educate the public on ways to prevent diseases and reduce health risks.
- **Advocacy and Education:** Support social justice through legal advocacy, social action, and supporting laws and measures that promote reform and protect civil rights, including election reform and tolerance among diverse groups.

- **Development and Relief Services:** Provide medical care and other human services as well as economic, educational, and agricultural development services to people around the world.
- **Humanitarian Relief Supplies:** Specialize in collecting donated medical, food, agriculture, and other supplies and distributing them overseas to those in need.
- **International Peace, Security, and Affairs:** Promote peace and security, cultural and student exchange programs, improve relations between particular countries, provide foreign policy research and advocacy, and United Nations-related organizations.
- **Religious Activities:** Support and promote various faiths.
- **Religious Media and Broadcasting:** Support organizations of all faiths that produce and distribute religious programming, literature, and other communications.
- **Non-Medical Science & Technology Research:** Support research and services in a variety of scientific disciplines, advancing knowledge and understanding of areas such as energy efficiency, environmental and trade policies, and agricultural sustainability.
- **Social and Public Policy Research:** Support economic and social issues impacting our country today, educate the public, and influence policy regarding healthcare, employment rights, taxation, and other civic ventures.

## D  Demographic Breakdown

Table 1 provides a detailed demographic breakdown per group.

## E  Detailed Qualitative Cognitive Load Breakdown

We provide additional details on the six cognitive dimensions. Among all dimensions, we also provide the codes representing different types of demand in a table form. The shaded cells represent the percentage of participants citing each source of mental demand, allowing for comparison within columns. The abbreviations in the columns: ST (Short Text Interface), S2P (Short Two-phase Interface), LT (Long Text Interface), and L2P (Long Two-phase Interface). Short and Long refer to the sum across both interfaces; Text and 2P (Two-phase interface) refer to the sum across both survey lengths. We include Sparklines for comparisons across these experiment groups. Future studies can use these as initial codebooks to conduct interface studies on preference construction.

### E.1  Sources of Mental Demand

Mental demand refers to the amount of mental and perceptual activity required to complete a task. Table 2 lists all qualitative codes and Figure 22 shows the boxplot of participant's subscale response. For thematic groups, we grouped them as source of demand (e.g., tracking remaining credits) and also of scope (e.g., Operational) as separated by the light gray line within each row.

### E.2  Sources of Physical Demand

Physical demand refers to the physical effort required to complete a task, such as physical exertion or movement. Most participants reported minimal physical demand ($N = 32$), reflected in the low NASA-TLX physical demand scores (Figure 23). Notably, 11 out of 20 participants who used the two-phase interface mentioned physical demand from using the mouse, reflecting interacting with

**Table 1: Participant Age and Gender Distribution by Experimental Condition**

| Condition | Mean Age | SD | Range | 25th | Median | 75th | Male | Female | Non-binary |
|---|---|---|---|---|---|---|---|---|---|
| Short Text | 31.6 | 13.7 | 18−67 | 23.8 | 29.5 | 32.8 | 4 | 6 | 0 |
| Short Two-Phase | 32.1 | 14.0 | 18−52 | 20.3 | 27.0 | 44.5 | 4 | 6 | 0 |
| Long Text | 36.0 | 14.8 | 21−61 | 24.0 | 33.0 | 42.8 | 2 | 7 | 1 |
| Long Two-Phase | 38.8 | 19.6 | 19−71 | 25.0 | 28.5 | 53.0 | 2 | 8 | 0 |



**Figure 22: Mental Demand Raw Score: Across all four experiment groups, participants' reported mental demand is spread across a wide range with many participants experiencing high mental demand.**

**Figure 23: Physical Demand Raw Score: Participants other than the long two-phase interface reported minimal physical demand. The long two-phase interface had the highest physical demand, likely due to increased mouse clicks and extended time spent looking at the vertical screen.**
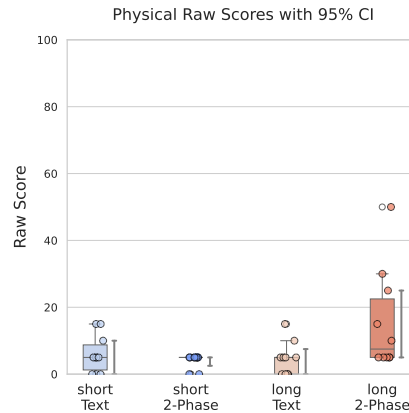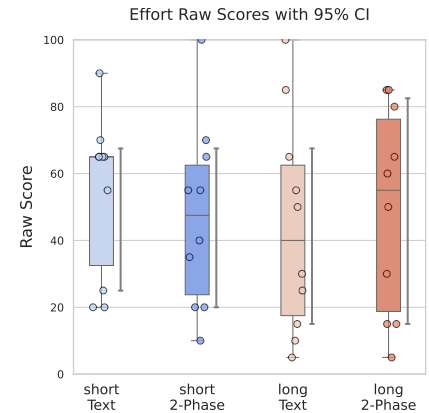
**Figure 24: Effort Raw Score: Effort scores show indifference across groups. All groups had high variance of responses indicating some participants requires high amount of effort when completing QS regardless of length and interface**

two interfaces. This is further supported by the raw NASA-TLX physical demand scores (Figure 23), which show a significant visual difference between short and long two-phase interfaces as well as between text and two-phase interfaces in long surveys. Table 3 presents all the relevant codes across experiment groups.

### E.3 Source of Effort

Effort refers to how hard participants felt they worked to achieve the level of performance they did. Since effort includes both mental and physical resource intensity, refer to Appendix E.1 and Appendix E.2 for definitions. Raw NASA-TLX effort scores (Figure 24) showed a similar spread across experiment groups, the qualitative analysis showed more distinction that participants using the two-phase interface considered options more comprehensively and felt less effort on completing operational tasks, similar to what we found on mental demands (Section E.1). For this subscale, we grouped codes through the lens of scope. Table 4 contains codes.

14 of the 20 participants using the text interface mentioned operational tasks as a source of effort, compared to 7 participants using the two-phase interface, with the lowest mention in the long two-phase interface group ($N = 2$).

*I wanted to bump up (an option) maybe to 4 or <option> to 5 and realize I couldn't. […] that would be effort came in of how do I want to really rearrange this to make it (the budget spending) maximize?*

💬 S029 (ST)

In contrast, strategic planning was reported as an effort source by 11 participants in the text interface, compared to 17 participants in the two-phase interface, with nearly all participants in the long two-phase interface group ($N = 9$) expressing effort related to it. In this subscale, we further categorize strategic planning into *narrow* and *broad* scopes as we did for mental demand (Appendix E.1). Participants using the two-phase interface ($N = 7$) had nearly mentioned double ($N = 4$) times regarding global strategies. For example:

*[…] the effort was how to rank order these (options) and allocate the resources behind the upvotes so that I can accurately depict what I want …say, a committee to focus on and allocate actual fungible resources, too.*

💬 S019 (L2P)

### E.4 Source from Performance

Performance refers to a person's perception of how successfully they have completed a task. Lower values indicate good perceived performance, while higher values suggest poor perceived performance. Raw NASA-TLX scores (Figure 25) show that participants had similar performance scores, although we highlighted nuanced differences in the main text. In addition to the differences mentioned in the main text, an interesting theme that emerged across experimental conditions was that participants' identified that *Social*

**Table 2: This table lists all the causes participants mentioned as contributing to their Mental Demand.**

| [ Mental Demand ] | Total | Version | | | | Experiment Conditions | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | ST | S2P | LT | L2P | Short | Long | Text | 2P |
| **Budget Management** | 14 | 3 | 3 | 5 | 3 | 6 | 8 | 8 | 6 |
| Budget within limited credit | 5 | 2 | 2 | 1 | 0 | 4 | 1 | 3 | 2 |
| Track remaining credits | 10 | 2 | 2 | 3 | 3 | 4 | 6 | 5 | 5 |
| Maximize credit usage | 8 | 2 | 3 | 2 | 1 | 5 | 3 | 4 | 4 |
| Operational | 12 | 3 | 2 | 4 | 3 | 5 | 7 | 7 | 5 |
| Strategic | 7 | 2 | 4 | 1 | 0 | 6 | 1 | 3 | 4 |
| **Preference Construction** | 39 | 10 | 9 | 10 | 10 | 19 | 20 | 20 | 19 |
| Determining relative preference | 16 | 4 | 4 | 5 | 3 | 8 | 8 | 9 | 7 |
| Option prioritization | 17 | 6 | 4 | 3 | 4 | 10 | 7 | 9 | 8 |
| Precise resource allocation | 30 | 9 | 6 | 9 | 6 | 15 | 15 | 18 | 12 |
| Narrow - Consider a few options/personal causes | 23 | 6 | 6 | 8 | 3 | 12 | 11 | 14 | 9 |
| Broad - Considering all options or higher order values | 23 | 5 | 5 | 4 | 9 | 10 | 13 | 9 | 14 |
| **Demand from Experiment Setup** | 24 | 6 | 6 | 6 | 6 | 12 | 12 | 12 | 12 |
| Many options on the survey | 6 | 0 | 0 | 3 | 3 | 0 | 6 | 3 | 3 |
| QS Mechanism | 4 | 2 | 0 | 2 | 0 | 2 | 2 | 4 | 0 |
| Recalling experience or understanding options | 20 | 5 | 6 | 4 | 5 | 11 | 9 | 9 | 11 |
| **Justification or Reflection on response** | 8 | 2 | 2 | 1 | 3 | 4 | 4 | 3 | 5 |
| **External Factors** | 12 | 3 | 1 | 4 | 4 | 4 | 8 | 7 | 5 |
| **Demand due to Interface** | 8 | 2 | 2 | 0 | 4 | 4 | 4 | 2 | 6 |
| Increase | 4 | 1 | 1 | 0 | 2 | 2 | 2 | 1 | 3 |
| Decrease | 4 | 1 | 1 | 0 | 2 | 2 | 2 | 1 | 3 |

**Table 3: Physical Demand Causes: Most participants expressed little or no physical demand. Results reflected that participants in the long two-phase interface required more actions, hence the higher mention of mouse usage as a source.**

| [ Physical ] | Total | Version | | | | Experiment Conditions | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | ST | S2P | LT | L2P | Short | Long | Text | 2P |
| **Reading** | 4 | 0 | 2 | 1 | 1 | 2 | 2 | 1 | 3 |
| **Mouse** | 16 | 3 | 5 | 2 | 6 | 8 | 8 | 5 | 11 |
| **Vertical Screen** | 4 | 1 | 0 | 1 | 2 | 1 | 3 | 2 | 2 |
| **None/Little** | 32 | 8 | 9 | 8 | 7 | 17 | 15 | 16 | 16 |

**Table 4: Effort Sources: Participants using the text interface focused more on operational tasks, while those using the two-phase interface focused more on strategic planning.**

| [ Effort ] | Total | Version | | | | Experiment Conditions | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | ST | S2P | LT | L2P | Short | Long | Text | 2P |
| **Operational** | 21 | 6 | 5 | 8 | 2 | 11 | 10 | 14 | 7 |
| **Strategic** | 28 | 6 | 8 | 5 | 9 | 14 | 14 | 11 | 17 |
| Narrow | 22 | 4 | 7 | 5 | 6 | 11 | 11 | 9 | 13 |
| Broad | 11 | 2 | 3 | 2 | 4 | 5 | 6 | 4 | 7 |
| **None/Little/a bit** | 9 | 2 | 1 | 3 | 3 | 3 | 6 | 5 | 4 |

**Table 5: Performance Causes: Most causes are shared across experiment conditions. We provided qualitative interpretations of their own performance assessments.**

| [ Performance ] | Total | Version | | | | Experiment Conditions | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | ST | S2P | LT | L2P | Short | Long | Text | 2P |
| **Operational Action** | 13 | 2 | 3 | 3 | 5 | 5 | 8 | 5 | 8 |
| Budget Control | 6 | 1 | 1 | 2 | 2 | 2 | 4 | 3 | 3 |
| Preference Reflection | 6 | 1 | 1 | 2 | 2 | 2 | 4 | 3 | 3 |
| Limited Resources | 5 | 1 | 2 | 1 | 1 | 3 | 2 | 2 | 3 |
| **Social Responsibility** | 8 | 2 | 2 | 2 | 2 | 4 | 4 | 4 | 4 |
| Decision maker | 7 | 1 | 2 | 2 | 2 | 3 | 4 | 3 | 4 |
| Outcome Uncertainty | 7 | 1 | 2 | 2 | 2 | 3 | 4 | 3 | 4 |
| **Performance Assessment** | | | | | | | | | |
| Did their best | 8 | 2 | 1 | 3 | 2 | 3 | 5 | 5 | 3 |
| Feel Good | 17 | 3 | 5 | 3 | 6 | 8 | 9 | 6 | 11 |
| Good Enough | 10 | 2 | 2 | 3 | 3 | 4 | 6 | 5 | 5 |

**Table 6: Temporal Demand Sources: Decision-making and Operational Tasks are the main causes. Participants framed their decision-making sources differently.**
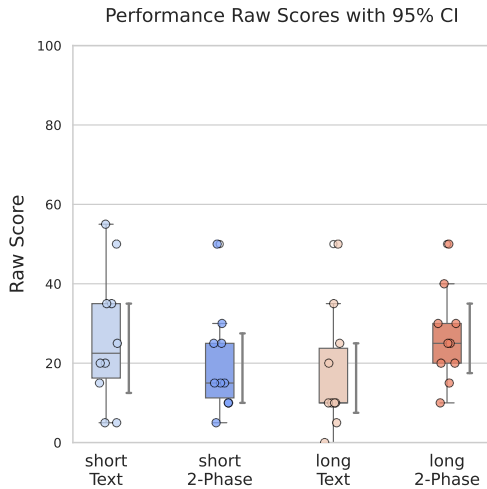
| [ Temporal ] | Total | Version | | | | Experiment Conditions | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | ST | S2P | LT | L2P | Short | Long | Text | 2P |
| **Budget Management** | 4 | 0 | 1 | 1 | 2 | 1 | 3 | 1 | 3 |
| **Decision Making** | 15 | 5 | 2 | 3 | 5 | 7 | 8 | 8 | 7 |
| Affirmative | 9 | 0 | 2 | 2 | 5 | 2 | 7 | 2 | 7 |
| Negative | 8 | 5 | 1 | 2 | 0 | 6 | 2 | 7 | 1 |
| **Operational** | 16 | 5 | 6 | 3 | 2 | 11 | 5 | 8 | 8 |
| Task completion | 8 | 2 | 2 | 3 | 1 | 4 | 4 | 5 | 3 |
| Being efficient | 8 | 3 | 4 | 0 | 1 | 7 | 1 | 3 | 5 |

**Table 7: Frustration Sources: Frustration comes from different levels of strategic operations or operational tasks.**

| [ Fustration ] | Total | Version | | | | Experiment Conditions | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | ST | S2P | LT | L2P | Short | Long | Text | 2P |
| **Strategic** | 17 | 4 | 4 | 5 | 4 | 8 | 9 | 9 | 8 |
| Higher-level | 11 | 3 | 2 | 3 | 3 | 5 | 6 | 6 | 5 |
| x Conflict between personal preference and broader society and common values | 6 | 1 | 1 | 2 | 2 | 2 | 4 | 3 | 3 |
| x Trade-offs among all options | 8 | 3 | 1 | 2 | 2 | 4 | 4 | 5 | 3 |
| Lower-Level | 10 | 3 | 3 | 2 | 2 | 6 | 4 | 5 | 5 |
| x Conflict between personal preference and | 4 | 1 | 2 | 0 | 1 | 3 | 1 | 1 | 3 |
| x Trade-offs among a few options | 8 | 2 | 2 | 2 | 2 | 4 | 4 | 4 | 4 |
| **Operational** | 15 | 4 | 5 | 2 | 4 | 9 | 6 | 6 | 9 |
| Credit management | 6 | 2 | 3 | 1 | 0 | 5 | 1 | 3 | 3 |
| Adhering to the Quadratic Mechanism | 5 | 2 | 1 | 1 | 1 | 3 | 2 | 3 | 2 |
| Deciding number of votes for an option | 4 | 2 | 0 | 0 | 2 | 2 | 2 | 2 | 2 |
| Making multiple decisions | 3 | 2 | 0 | 0 | 1 | 2 | 1 | 2 | 1 |
| Understanding Option | 4 | 0 | 3 | 0 | 1 | 3 | 1 | 0 | 4 |
| **None/Little** | 16 | 4 | 5 | 5 | 2 | 9 | 7 | 9 | 7 |

*Responsibility* influenced their performance scores. Table 5 presents a detailed breakdown of our codes.

*Social Responsibility.* This theme refers to concerns about performance when participants reflected on how their final vote counts would be perceived by others ( S041 💬 *I don't want people to think that I just don't care about <ethnicity> people at all* ) or how their votes might influence real-world decision-making ( S027 💬 *Some of these things might . . . have outcomes that I didn't foresee* ).

Performance Raw Scores with 95% CI



**Figure 25: Performance Demand Raw Score: Participants showed indifferent performance raw scores across experiment conditions, all trending toward satisfactory.**

## E.5 Temporal Demand

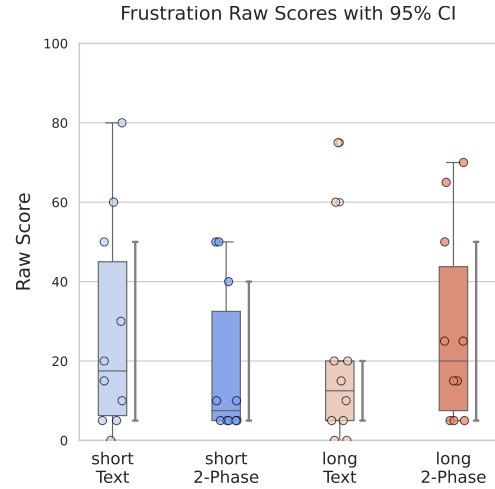Table 6 lists all the temporal demand codes.

## E.6 Frustration

Table 7 lists all codes related to participants' sources of frustration.

## F Additional voting behavior data

In this section, we describe additional voting behavior that we observed. The reason why we decided to focus on the percentage of remaining credits comes from prior literature 'scarcity frames value' [79], a driver that makes researchers believe makes quadratic voting more accurate [7]. We did not follow Quarfoot et al. [68] in counting accumulated votes over time due to varying total times across individuals.

We observed the number of vote adjustments given a remaining vote credit percentage. Figure 27 showed all the voting actions over the remaining credit for the four experiment conditions. Here we see two distinct patterns between the short survey and the long survey in terms of participant behaviors. In long surveys, participants exhibited more actions both when the budget was abundant and when it began to run out. This pattern was more pronounced with the long two-phase interface. This difference is why we further focused on the long QS group.

Frustration Raw Scores with 95% CI



**Figure 26: Frustration Raw Score: Participants other than the long text interface highlighted several operational tasks that led to frustration. All groups share causes from strategic planning.**

Figure 28 presents the comparison between when participants make small or large vote adjustments at different budget levels. Revisiting the KDE curve in the second row in Figure 27 and the curve of the second row in Figure 28 show a stronger bimodal distribution for small vote adjustments across interfaces. In fact, the bimodal distribution is more pronounced in the two-phase interface. This suggests that participants make small adjustments both at the beginning and toward the end of the QS. However, the two-phase interface shows more frequent and faster edits towards the end. In comparison, participants also made more large vote adjustments early on that spread more equally compared to the text interface. This indicates that participants had a clearer idea of how to distribute their credits across the options.

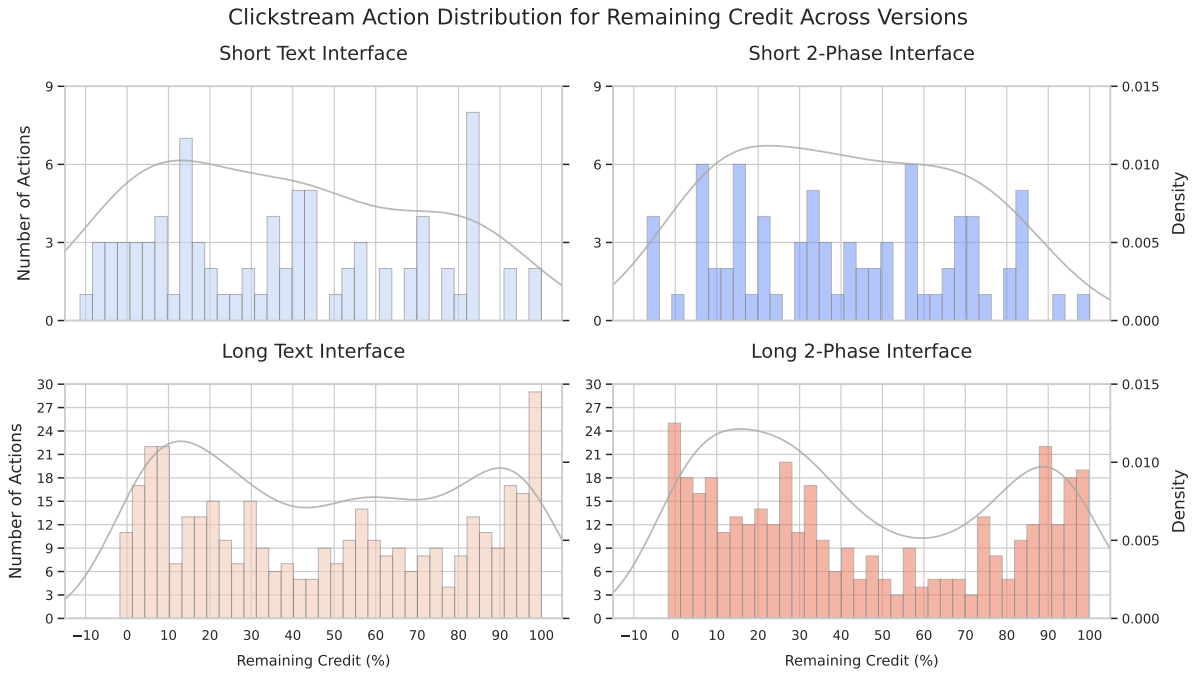## G Modeling NASA-TLX Weighted Scores and Subscales

This section first describes the hierarchical Bayesian ordinal regression model used for the NASA-TLX weighted scores and subscales. We then present the results for each subscale.
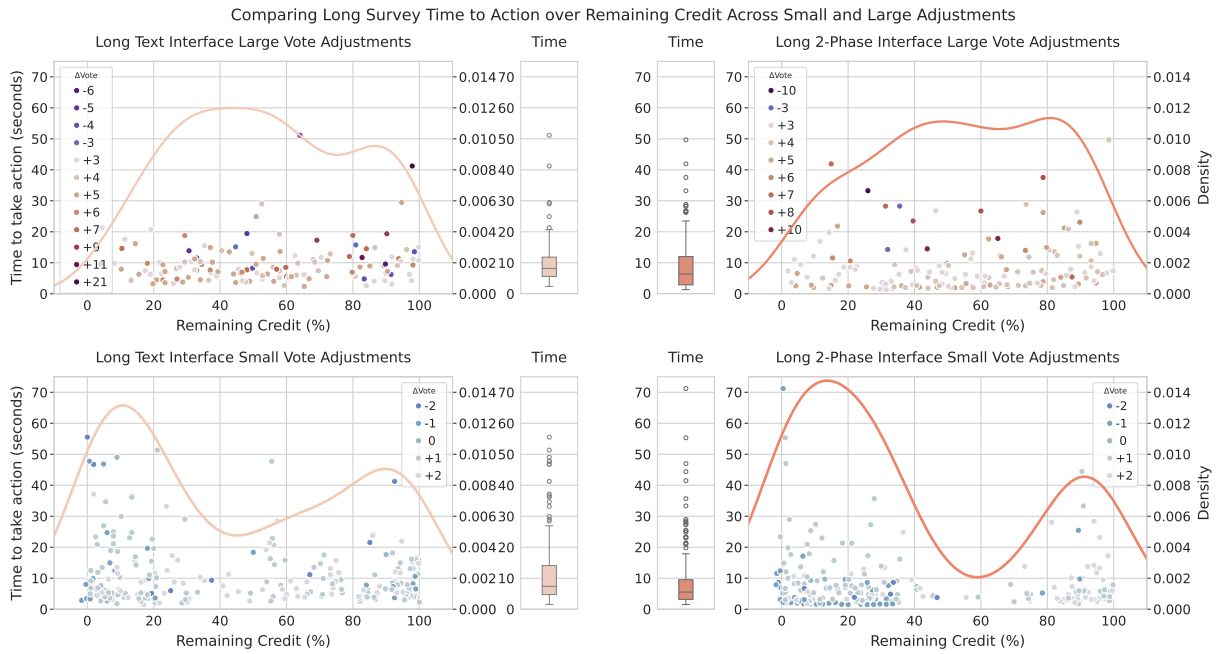
## G.1 Modeling Approach

### G.1.1 Dependent variables.

*NASA-TLX weighted scores.* are transformed from a continuous 0–100 scale to cognitive levels: low, medium, somewhat high, high, and very high, as described by Hart and Staveland [34]. This transformation helps the model adapt to sparse data. In our study, there were no participants who expressed "low" or "very high"; thus, we modeled the predictive variables as "medium," "somewhat high," and "high."

*NASA-TLX subscale ratings.* are transformed into ordinal groups using minimum frequency binning [25]. Minimum frequency binning involves grouping adjacent response categories until each bin

**Figure 27: This plot counts the number of voting actions when there are $x$ percentages of credits remaining. A KDE plot is provided to help better understand the action distribution.**



**Figure 28: This plot further separates participants' interaction behavior based on the number of votes participants adjusted. We observed a bimodal interaction pattern across long QS when small vote adjustments are made.**

meets a predefined minimum number of observations. Since the subscale uses a 21-point Likert scale and we have 40 participants, the data are very sparse. Minimum frequency binning mitigates this by ensuring similar numbers of participants in each bin. We applied weighted bins across all participants within the same subscale, ensuring that each bin contained at least 10 participants.

*Likelihood.* With these ordinal outcome variables, we designed $y_i$ as the observed ordinal category for participant $i$. Then:

$$y_i \sim \text{OrderedLogistic}(\eta_i, \boldsymbol{\tau}), \tag{1}$$

where $\eta_i$ is the latent predictor, and $\boldsymbol{\tau}$ denotes the cutpoints demarcating the boundaries between the ordinal categories as in Equation (2). The cutpoints $\boldsymbol{\tau}$ ensure that $\tau_1 < \tau_2 < \cdots < \tau_{K-1}$ by construction.

$$\boldsymbol{\tau} \sim \text{OrderedTransform}(\mathcal{N}(0,1)^{K-1}), \tag{2}$$

*G.1.2 Independent Variables and latent predictor.* For this model, we used three independent variables: length ($\gamma_i$, an ordinal variable), interface type ($\beta_I$, a categorical variable), and the interaction between the two ($\phi_{i,j}$) to construct the latent predictor $\eta_i$. Specifically, the latent predictor $\eta_i$ is constructed as:

$$\eta_i = \alpha + \gamma_i + \beta_I[I_i] + \phi_{i,j}, \tag{3}$$

where: $\alpha$ is a global intercept drawn from $\mathcal{N}(0,1)$, $\gamma_i$ captures the (ordinal) effect of length, $\beta_I[I_i]$ is the effect for interface $I_i$, and $\phi_{i,j}$ is the interaction between length $i$ and interface $j$.

Since length has two levels (short and long), we define the following equation to account for ordinality:

$$\gamma_i = \mu_L + \beta_L \cdot L_i \tag{4}$$

where $L_i \in \{0, 1\}$, making $\gamma_i = \mu_L$ for the short condition and $\gamma_i = \mu_L + \beta_L$ for the long condition. We assign standard normal priors to these parameters: $\mu_L \sim \mathcal{N}(0,1)$ and $\beta_L \sim \mathcal{N}(0,1)$.

*Interface Effects.* We model the interface effects using a non-centered parameterization to improve numerical stability and encourage partial pooling across the two interface levels. Specifically, we let $\mu_{\beta_I} \sim \mathcal{N}(0,1)$ and $\sigma_{\beta_I} \sim \text{Exponential}(1)$ represent the shared mean and scale of the interface effects. We then sample a raw effect vector $\beta_{I_{\text{raw}}} \sim \mathcal{N}(0,1)^2$. Combining these, we define:

$$\beta_I = \mu_{\beta_I} + \sigma_{\beta_I} \cdot \beta_{I_{\text{raw}}} \tag{5}$$

where $\beta_I \in \mathbb{R}^2$ contains the effect for each of the two interface levels, and $\beta_I[I_i]$ indexes the effect for participant $i$'s interface.

*Interaction Effects.* To capture potential interaction effects between length and interface types, we assign one interaction parameter, $\phi_{i,j}$, to each combination of length $i$ and interface $j$. Rather than sampling these $\phi_{i,j}$ directly, we employ a non-centered parameterization:

$$\boldsymbol{\phi} = L_\Omega \left( \sigma_\phi \odot z_\phi \right),$$

where $\boldsymbol{\phi}$ is a $2 \times 2$ matrix of interaction parameters (since we have 2 levels of length and 2 levels of interface), $z_\phi \sim \mathcal{N}(0,1)^{2\times2}$, $\sigma_\phi \sim \text{Exponential}(1)^{2\times2}$, and $L_\Omega$ is the Cholesky factor of a correlation matrix drawn from an LKJ(2) prior. We then define

$$\phi_{i,j} = \left[ \boldsymbol{\phi} \right]_{i,j},$$

making $\phi_{i,j}$ a *single scalar* drawn from the correlated matrix $\boldsymbol{\phi}$.

*G.1.3 Posterior predictive plots.* We conducted the Bayesian analysis using NumPyro, a widely used framework for Bayesian inference. We used Markov Chain Monte Carlo (MCMC) sampling, a method commonly applied in Bayesian inference. The model converged successfully, as evidenced by an $\hat{R}$ value of 1 for each subscale and the overall weighted TLX scores, indicating that multiple sampling chains converged. We plotted the posterior predictive distribution of the model to compare the observed data with the model's predictions. Figure 29 shows the posterior predictions vs. observed data for the six subscales.

## G.2 Model Results

*G.2.1 Mental Subscale.* Figure 30 shows pairwise Bayesian results from mental demand highlighted 70.4% of posterior probability that participants in the long two-phase condition had a higher mental demand compared to the short two-phase condition. On the other hand, the short text condition had a 74.5% posterior probability of having a higher mental demand compared to the short two-phase condition. This is additional evidence that prompted us to believe that the participants in the short two-phase participants benefited from the organization phase. The sheer number of added options in the long two-phase condition may have added additional demand to participants, leading to higher mental demand.
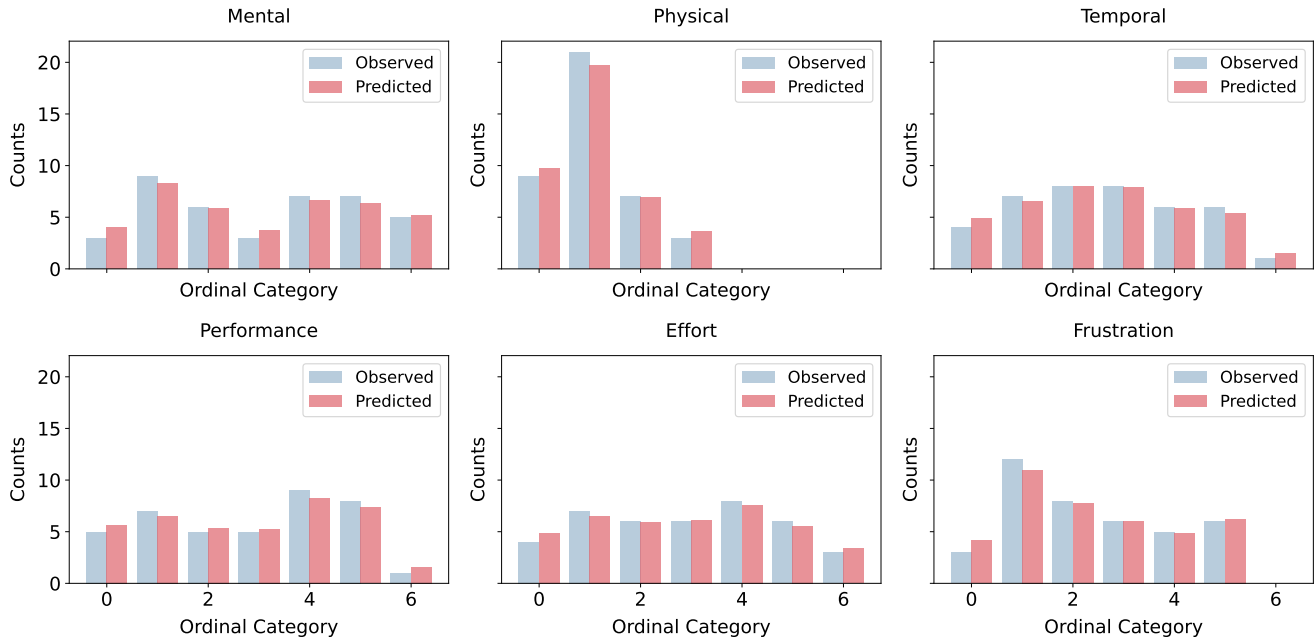
*G.2.2 Physical Subscale.* Figure 31 shows the pairwise comparison of the physical subscale. Notable results show that there is a 86.1% posterior probability that the long text condition had a lesser physical demand compared to the short text condition. This is counter intuitive as the long text participants actually traversed much higher edit distances. We are not clear what prompted their self reported value and requires future research.

*G.2.3 Temporal Subscale.* Figure 32 shows the pairwise comparison of the temporal subscale. The results show that the long two-phase condition once again had a 74.6% posterior probability of having a lower temporal demand compared to the short text condition. Conversely, participants in the long two-phase condition had a 71.1% posterior probability of having a higher temporal demand compared to the short two-phase condition, reflecting the longer time they took to complete the survey questions. We believe that the lower temporal demand in the long two-phase condition is potential indicator of the participants' satisficing behavior.
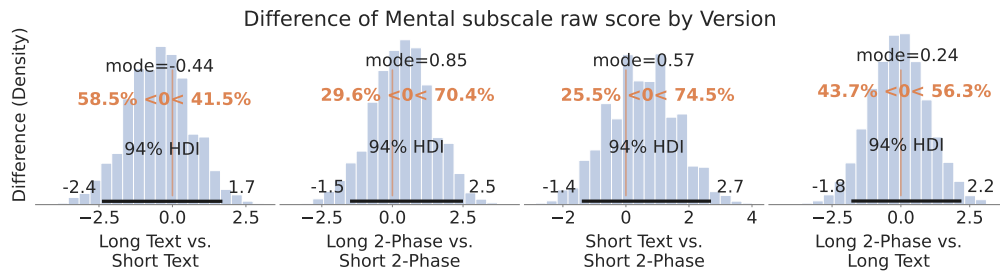
*G.2.4 Performance Subscale.* We omit the pairwise comparison of the performance subscale due to the mixed signals. We focused on the qualitative results analyzed in the main text.

*G.2.5 Effort Subscale.* We omit the pairwise comparison of the effort subscale due to its similarity to the mental demand subscale.
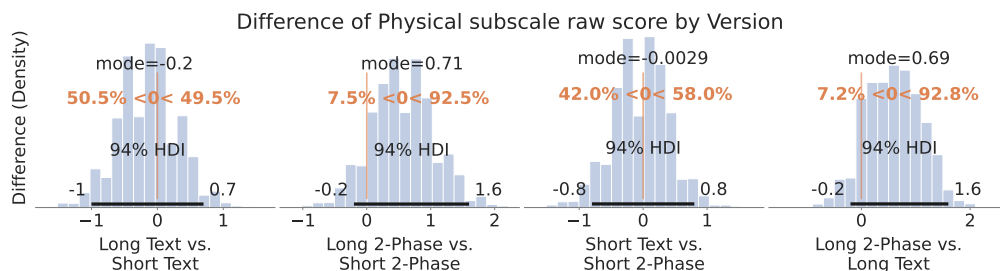
*G.2.6 Frustration Subscale.* Figure 33 shows the pairwise comparison of the frustration subscale. The results show that the long two-phase condition had a 68.3% posterior probability of having a higher frustration compared to the short two-phase condition, likely due to the added number of options to assess.
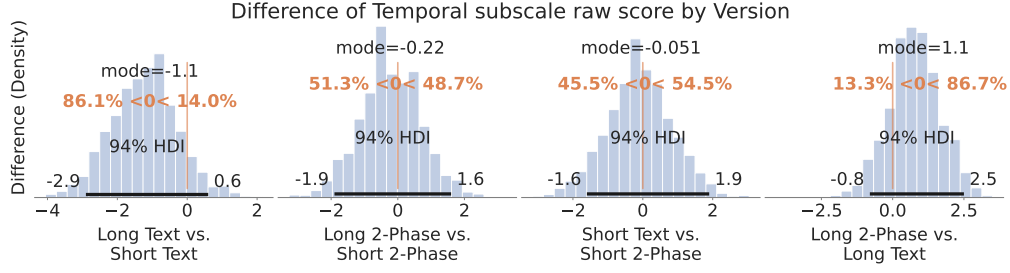
Figure 29: Posterior Predictions vs. observed data for NASA-TLX subscales. The plot shows the observed number of participants in each bin compared to the posterior predictions from the model. Takeaway of the plot: We believe that the model is reasonable at capturing the distribution of the subscales given the sparsity of the data.
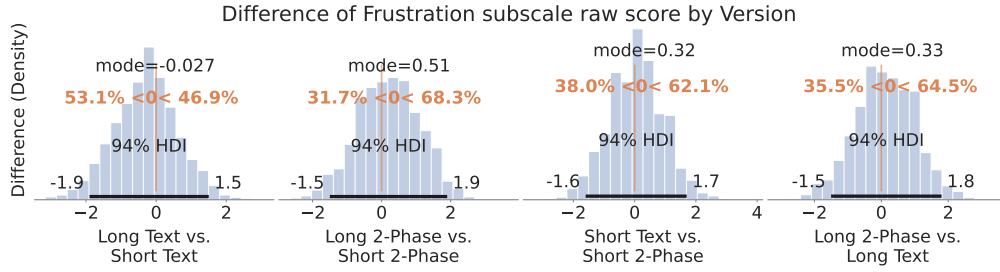


Figure 30: Differences in the mental subscale scores by version. Main Takeaway: Participants in the long two-phase condition show trends to increase mental demand compared to the short two-phase. Within the short text condition, participants in the short two-phase condition show a trend to reduce mental demand.



Figure 31: Differences in the physical subscale scores by version. Main Takeaway: Participants in the long two-phase condition show trends to increase physical demand compared to short two-phase and long text despite the long text participants traversing higher edit distances.

**Figure 32: Differences in the temporal subscale scores by version. Main Takeaway: Participants in the long text condition show a trend that it reduces temporal demand compared to the short text condition and the long two-phase condition.**



**Figure 33: Differences in the frustration subscale scores by version. Main Takeaway: The model does not see a significant difference in the frustration subscale between experiment groups other than a trend for participants in the long two-phase condition to have higher frustration than the short two-phase participants.**

## H Modeling Total Time

### H.1 Dependent Variables

The dependent variable is the total time $T_i$ spent on option $i$ measured in seconds. This measure captures both the duration participants took to vote and, where applicable, the time they spent organizing or reordering their options beforehand. We categorize the data into four experimental conditions: Short Text, Short Two-Phase, Long Text, and Long Two-Phase. These conditions are indexed by $k$, fit using separate submodels.

### H.2 Modeling Approach

We modeled the total time for each experimental condition using separate Gamma likelihood models. The Gamma distribution is well-suited for modeling positive continuous data, such as time measurements, which are often skewed and strictly positive. Equation 6 shows the model for the total time. The shape parameter $\alpha_k$ and rate parameter $\beta_k$ were each assigned priors drawn from their own Gamma distributions, as described in Equations 7 and 8.

$$T_i \sim \text{Gamma}(\alpha_k, \beta_k) \tag{6}$$

$$\alpha_k \sim \text{Gamma}(2.0, 0.5) \tag{7}$$

$$\beta_k \sim \text{Gamma}(1.0, 1.0) \tag{8}$$

## I Modeling Edit Distance

This section presents our hierarchical Bayesian approaches for analyzing the edit distance data. We first describe a model for edit

distance per option (Appendix I.1), followed by analysis for edit distance per action (Appendix I.2). Finally, we detail a model for cumulative edit distances (Appendix I.3).

### I.1 Model 1: Edit Distance per Option

*I.1.1 Likelihood.* The dependent variable in this model is the edit distance accumulated for each option, denoted by $D_i$, where $i$ refers to the $i$-th observation. Since $D_i$ must be positive, we model it using an exponential likelihood:

$$D_i \sim \text{Exponential}(\text{scale} = \lambda_i). \tag{9}$$

*I.1.2 Independent variables and regression model.* We designed $\eta_i$ as the linear predictor that informs $D_i$ through the following transformation:

$$\lambda_i = \exp(\eta_i), \tag{10}$$

where $\lambda_i$ is the scale (i.e., mean) parameter of the Exponential distribution, and thus must be positive.

This linear predictor:

$$\eta_i = \gamma_i + \beta_I[I_i] + \phi_{ij} + U_i \tag{11}$$

consists of four components: the length of the option $L_i$, interface type $I_i$, and interaction effect between both length and interface $\phi_{ij}$, and user effect $U_i$ which we describe in the following paragraphs.

*Length.* Since length has two levels (short and long), we define:

$$\gamma_i = \mu_L + \beta_L \cdot L_i \tag{12}$$

where $L_i \in \{0, 1\}$, making $\gamma_i = \mu_L$ for the short condition and $\gamma_i = \mu_L + \beta_L$ for the long condition. We assign standard normal priors to these parameters: $\mu_L \sim \mathcal{N}(0, 1)$ and $\beta_L \sim \mathcal{N}(0, 1)$.

*Interface.* We model the interface effects using a non-centered parameterization to improve numerical stability and encourage partial pooling across the two interface levels. Specifically we let $\mu_{\beta_I} \sim \mathcal{N}(0, 1)$ and $\sigma_{\beta_I} \sim \text{HalfNormal}(0.5)$ represent the shared mean and scale of the interface effects. We then sample a raw effect vector $\beta_{I_{\text{raw}}} \sim \mathcal{N}(0, 1)^2$. Combining these, we define:

$$\beta_I = \mu_{\beta_I} + \sigma_{\beta_I} \cdot \beta_{I_{\text{raw}}} \tag{13}$$

where $\beta_I \in \mathbb{R}^2$ contains the effect for each of the two interface levels, and $\beta_I[I_i]$ indexes the effect for participant $i$'s interface.

*Interaction Effects.* To capture potential interaction effects between length and interface types, we assign one interaction parameter, $\phi_{i,j}$, to each combination of length $i$ ($i \in \{0, 1\}$) for short and long surveys and interface $j$ ($j \in \{0, 1\}$) for the two interface types. Rather than sampling these $\phi_{i,j}$ directly, we employ a non-centered parameterization:

$$\boldsymbol{\phi} = L_{\Omega} \left( \sigma_{\phi} \odot z_{\phi} \right),$$

where $\boldsymbol{\phi}$ is a $2 \times 2$ matrix of interaction parameters (since we have 2 levels of length and 2 levels of interface), $z_{\phi} \sim \mathcal{N}(0, 1)^{2 \times 2}$, $\sigma_{\phi} \sim \text{HalfNormal}(0.5)^{2 \times 2}$, and $L_{\Omega}$ is the Cholesky factor of a $2 \times 2$ correlation matrix drawn from an $\text{LKJ}(2)$ prior with shape parameter $\eta = 3$. We then define

$$\phi_{ij} = \left[ \boldsymbol{\phi} \right]_{i,j} \tag{14}$$

making $\phi_{ij}$ a *single scalar* drawn from the correlated matrix $\boldsymbol{\phi}$.

*Individual user effects.* Similar to the interface, we also applied a non-centered parameterization to user effects using the same approach:

$$U_i = \mu_U + \sigma_U \cdot z_U \tag{15}$$

We assign weakly informative priors for the user effects: $\mu_U \sim \mathcal{N}(0, 1)$ and $\sigma_U \sim \text{Exponential}(0.5)$, which represent the shared mean and scale of the user effects. We use $z_U \sim \mathcal{N}(0, 1)^{40}$. to denote the 40 participant's raw user effect vector. This approach allow us to capture user variations across all users.
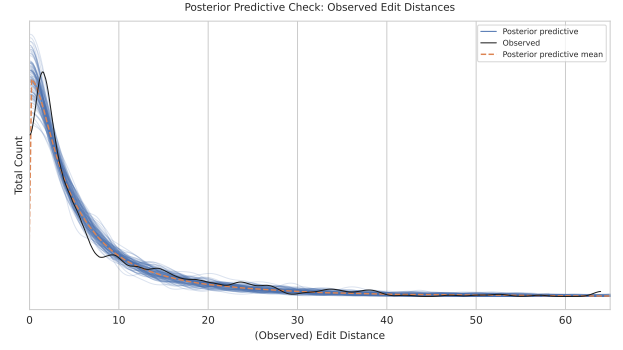
*I.1.3 Posterior predictive plots.* Our Bayesian model converged successfully, as evidenced by an $\hat{R}$ value of 1 in the model summary. We plotted the posterior predictive distribution for the edit distance per option in Figure 34. This figure compares the models posterior predictive distribution with the observed data.

## I.2 Model 2: Edit Distance with Separate Mean and Variance Predictors

*I.2.1 Likelihood.* The dependent variable for this model is the edit distance $D_i$, where positive values indicate a downward movement and negative values indicate an upward movement. To allow for different effects on both the mean and variance, we model $D_i$ using a Normal distribution:

$$D_i \sim \mathcal{N}\left( \mu_i, \sigma_{\text{obs},i} \right) \tag{16}$$

Because our aim is to capture potential differences in variability (e.g., hypothesizing that a two-phase interface might yield lower



Figure 34: Posterior Predictions vs. observed data for edit distance per option. Each blue line represents a draw from the posterior distribution, while the black line represents the observed data. Dotted line represents the mean of the posterior data. Takeaway of the plot: We believe that the model is reasonable at capturing the distribution.

oscillation than a text-based interface), we separately model both the mean $\mu_i$ and the standard deviation $\sigma_{\text{obs},i}$.

*I.2.2 Independent variables and regression model.* We specify two linear predictors: one for the mean $\mu_i$ (Equation 17) and one for the (logged) standard deviation $\log(\sigma_{\text{obs},i})$ (Equation 18). Both linear predictors incorporate the following factors: the length of the option $L_i$, the interface type $I_i$, an interaction term $\phi_{ij}$, and a user-specific term $U_i$.

$$\mu_i = \gamma_{\mu,i} + \beta_{I,\mu}[I_i] + \phi_{\mu,ij} + U_{\mu,i}, \tag{17}$$
$$\log(\sigma_{\text{obs},i}) = \gamma_{\sigma,i} + \beta_{I,\sigma}[I_i] + \phi_{\sigma,ij} + U_{\sigma,i}. \tag{18}$$

*Length ($L_i$).* Similar to the previous model, we continue to define length as an ordinal value. In this model, the effect for mean and variance are modeled separately. We write:

$$\gamma_{\mu,i} = \mu_{L,\mu} + \beta_{L,\mu} \cdot L_i, \tag{19}$$
$$\gamma_{\sigma,i} = \mu_{L,\sigma} + \beta_{L,\sigma} \cdot L_i. \tag{20}$$

For both the mean and variance parts, $\mu_{L,\mu}, \beta_{L,\mu}$ and $\mu_{L,\sigma}, \beta_{L,\sigma}$ capture how option length shifts the location and scale of the distribution, respectively. We assign weakly informative normal priors:

$$\mu_{L,\mu}, \beta_{L,\mu}, \mu_{L,\sigma}, \beta_{L,\sigma} \sim \mathcal{N}(0, 1). \tag{21}$$

*Interface ($I_i$).* We treat the interface type as a categorical variable with two levels. As in Model 1, we use a non-centered parameterization for numerical stability and partial pooling. For the mean part, we define:

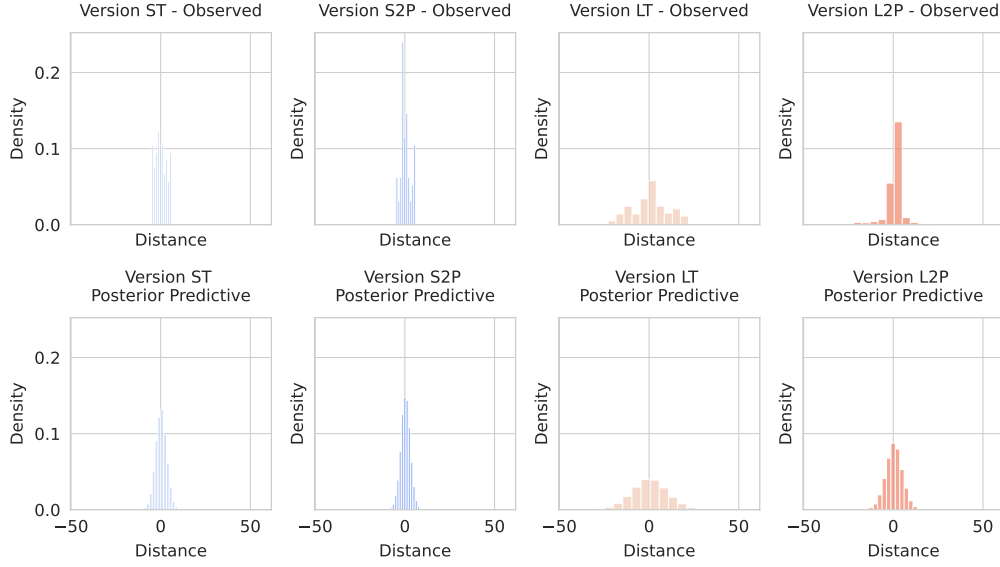$$\beta_{I,\mu}[I_i] = \mu_{I,\mu} + \sigma_{I,\mu} \cdot z_{I,\mu}[I_i]. \tag{22}$$

and similarly for the variance part:

$$\beta_{I,\sigma}[I_i] = \mu_{I,\sigma} + \sigma_{I,\sigma} \cdot z_{I,\sigma}[I_i]. \tag{23}$$

We place weakly informative priors on the intercepts:

$$\mu_{I,\mu}, \beta_{I,\mu}, z_{I,\mu}, \mu_{I,\sigma}, \beta_{I,\sigma}, z_{I,\sigma} \sim \mathcal{N}(0, 1), \tag{24}$$
$$\sigma_{I,\mu}, \sigma_{I,\sigma} \sim \text{HalfNormal}(0.5). \tag{25}$$

**Figure 35: Posterior Predictions vs. Observed Data for Edit Distance per Option. The first row represents the distribution of edit distance per version. The second row shows the posterior predictions after multiple draws Takeaway of the plot: We believe that the model is reasonable at capturing the shape of the distributions though being slightly conservative for extreme values at the center. Future model enhancements could re-model them with a student-t distribution.**

.

*Interaction Effects ($\phi_{ij}$).* We hypothesize that the effect of length might vary by interface. Similar to Model 1's approach, we employ a non-centered parameterization with an LKJ correlation prior. Specifically, for both the mean and variance parts, we define:

$$\phi_{\mu,ij} = \left[ L_{\Omega,\mu} \left( \sigma_{\phi,\mu} \odot z_{\phi,\mu} \right) \right] i, j, \tag{26}$$

$$phi\sigma, ij = \left[ L_{\Omega,\sigma} \left( \sigma_{\phi,\sigma} \odot z_{\phi,\sigma} \right) \right] i, j, \tag{27}$$

where $i \in 0, 1$ (short or long) and $j \in 0, 1$ (two interface types). We continue the use of weakly informed priors:

$$z_{\phi,\mu}, z_{\phi,\sigma} \sim \mathcal{N}(0,1), \sigma_{\phi,\mu}, \sigma_{\phi,\sigma} \sim \text{HalfNormal}(0.5), \tag{28}$$

$$L_{\Omega,\mu}, L_{\Omega,\sigma} \sim \text{LKJ}(3). \tag{29}$$

*Individual user effects ($U_i$).* To account for participant-level variability, we follow model 1 and adopt a non-centered parameterization but allow each user to have a distinct shift on both $\mu_i$ and $\log(\sigma_{\text{obs},i})$:

$$U_{\mu,i} = \mu_{U,\mu} + \sigma_{U,\mu} \cdot z_{U,\mu,i}, \tag{30}$$

$$U_{\sigma,i} = \mu_{U,\sigma} + \sigma_{U,\sigma} \cdot z_{U,\sigma,i}, \tag{31}$$

with priors:

$$\mu_{U,\mu}, \beta_{U,\mu}, z_{U,\mu,i}, \mu_{U,\sigma}, \beta_{U,\sigma}, z_{U,\sigma,i} \sim \mathcal{N}(0,1), \tag{32}$$

$$\sigma_{U,\mu}, \sigma_{U,\sigma} \sim \text{HalfNormal}(0.5). \tag{33}$$

*I.2.3 Posterior predictive plots.* Our Bayesian model converged successfully, as evidenced by an $\hat{R}$ value of 1 in the model summary. We plotted the posterior predictive distribution for the edit distance per option in Figure 35. This figure compares the models posterior predictive distribution with the observed data.

*I.2.4 Model Results.* Figure 36 shows the pairwise comparison of the variance of edit distance in the first row followed by the effect size in the second row. In addition to the comparison within the same survey length, we provide all pairwise comparisons. A notable result that we omit from the main text is that if we compare the variance between the long and short text, and the variance between the long and short two-phase, we see that the text group had three times the standard deviation compared to the two-phase group. This indicates that the organization phase minimize the edit distance despite the increase in survey length.

## I.3 Model 3: Long QS Cumulative Edit Distance

The dependent variable for this model is the cumulative edit distance $D_i$, a positive continuous variable measured at each step within a version for each user. We modeled this to test our hypothesis that for each participant, the growth rate of the edit distance is consistent. To accommodate its positive nature, we model $D_i$ using a Truncated Normal distribution:

$$D_i \sim \text{TruncatedNormal}(\mu_i, \sigma_{\text{obs},i}, \text{lower} = 0), \tag{34}$$

where the observation-specific standard deviation prior is:

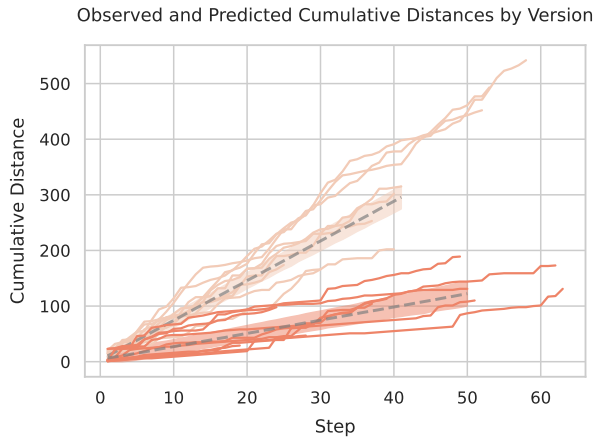$$\sigma_{\text{obs},i} \sim \text{HalfNormal}(0.3). \tag{35}$$

*I.3.1 Independent Variables and Regression Model.* We incorporate the following independent variables: the step number when completing QS ($S_i$), the interface version ($V_i$), and user-specific effects ($U_i$). The interface version and user-specific effects are modeled using hyperpriors to capture variability across groups.

The linear predictor for $D_i$ is given by:

$$\mu_i = \alpha_{\text{shared}} + \beta_v[V_i] \cdot S_i + U_i \cdot S_i, \tag{36}$$

Difference and Effect Size of Distance per Option by Version



Figure 36: Differences in the variance of edit distance by version. Main takeaway: This plot shows that with two-phase interface, there is a reduction in edit distance variance when the number of option grows.



Figure 37: Posterior Predictions vs. observed data for cumulative edit distance. The plot showed each observed user's cumulative edit distance in different shades for the two groups of participants. Dotted line represent the posterior predictive mean. Takeaway of the plot: We believe that the model is reasonable at capturing slop of the cumulative trends.

where $\alpha_{\mathrm{shared}}$ represents the global intercept, $\beta_v[V_i]$ models the interface version effects, and $U_i$ captures individual user-specific effects. The intercept is assigned the following prior:

$$\alpha_{\mathrm{shared}} \sim \mathcal{N}(2.0, 0.5). \tag{37}$$

*Interface Version ($V_i$).* Interface effects are modeled as:

$$\beta_v[V_i] \sim \mathcal{N}(\mu_\beta, \sigma_\beta), \tag{38}$$

where the hyperparameters for the interface effect distribution are:

$$\mu_\beta \sim \mathcal{N}(0.05, 0.05), \quad \sigma_\beta \sim \mathrm{HalfNormal}(0.1). \tag{39}$$

*User Effects ($U_i$).* Instead of directly sampling $U_i$, we follow the reparameterization approach:

$$U_i = \mu_U + \sigma_U \cdot z_{U,i}, \tag{40}$$

where we assign weakly informative priors $\mu_U \sim \mathcal{N}(0, 1)$ and $\sigma_U \sim \mathrm{HalfNormal}(0.1)$ to represent the shared mean and scale of the user effects. The term $z_{U,i} \sim \mathcal{N}(0, 1)$ captures individual user variability, allowing us to model deviations across users while maintaining a structured prior.

*I.3.2 Posterior Predictive Plots.* Our Bayesian model converged successfully, as indicated by an $\hat{R}$ value of 1 in the model summary. Figure 37 presents the posterior predictive distribution for cumulative edit distance, demonstrating alignment between the predicted and observed data.