

Addressing Repetition in Crowdsourcing: A Concept for Fast-Form Entry

van Berkel, Niels; Schneiders, Eike; Jacobsen, Rune Møberg

Published in:

Adjunct Proceedings of the ACM SIGCHI Conference on Human Factors in Computing Systems (CHI'22 EA)

Publication date:

2022

Document Version

Accepted author manuscript, peer reviewed version

[Link to publication from Aalborg University](#)

Citation for published version (APA):

van Berkel, N., Schneiders, E., & Jacobsen, R. M. (2022). Addressing Repetition in Crowdsourcing: A Concept for Fast-Form Entry. In Adjunct Proceedings of the ACM SIGCHI Conference on Human Factors in Computing Systems (CHI'22 EA) Association for Computing Machinery (ACM).

General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal -

Take down policy

If you believe that this document breaches copyright please contact us at vbn@aub.aau.dk providing details, and we will remove access to the work immediately and investigate your claim.

Addressing Repetition in Crowdsourcing: A Concept for Fast-Form Entry

NIELS VAN BERKEL, Aalborg University, Denmark

EIKE SCHNEIDERS, Aalborg University, Denmark

RUNE MØBERG JACOBSEN, Aalborg University, Denmark

This workshop paper outlines a conceptual browser plugin that enables crowdworkers to store and later rapidly provide personal information frequently requested in crowdsourcing tasks. Personal data, including demographic data such as age and ethnicity, as well as responses to commonly used personality-related survey instruments, is often critical to collect in crowdsourcing tasks but results in a repetitive experience for crowdworkers. From a requesters perspective, this repetition can result in reduced data quality or the decision to abstain from collecting extensive information on the workers completing a given task. Moreover, given the extensive role of crowdworkers in labelling training data for artificial intelligence applications, ensuring awareness of the workers' characteristics can help alleviate future biases. In this work, we present the motivation and design requirements for this (hypothetical) plugin and seek input from the community towards its future development.

CCS Concepts: • **Human-centered computing** → **Human computer interaction (HCI)**; *Field studies*.

Additional Key Words and Phrases: Crowdsourcing, fast-fill, autocomplete, survey, demographics, questionnaire

ACM Reference Format:

Niels van Berkel, Eike Schneiders, and Rune Moberg Jacobsen. 2018. Addressing Repetition in Crowdsourcing: A Concept for Fast-Form Entry. In *CHI REGROW '22: Reimagining Global Crowdsourcing for Better Human-AI Collaboration, April 30–May 6, 2022, New Orleans, USA*. ACM, New York, NY, USA, 6 pages. <https://doi.org/10.1145/1122445.1122456>

1 INTRODUCTION

Crowdworkers play an essential role in today's data-labelling and scientific data collection infrastructures. The extensive datasets that power Artificial Intelligence (AI) systems typically require substantial amounts of manual labelling. Similarly, researchers rely on human data input to test their hypotheses. This staggering increase in the amount of manual human input required has led to the uptake of various crowdworker markets, such as the well known Amazon Mechanical Turk, Prolific, Toloka, and others. A recent investigation into the online labour market highlights the massive nature of these operations – involving tens of millions of individuals, and demand for their services growing at 20% year per year [11].

Amidst this growing reliance on crowdworkers, there is an increasing consideration for, and realisation of malpractices concerning, the working conditions of these individuals [3, 7, 10, 15]. This includes, for example, critique on the financial compensation of crowdworkers [7, 15], the sometimes explicit or extreme content with which they are presented [5, 8], and how crowdsourcing platforms are not allowing for time-flexible completion of tasks [12]. In addition to critique on the worker conditions in crowdsourcing, the use of this data collection method has also received scrutiny from within the research community – in particular pointing to the biases that may emerge in the resulting AI systems when trained on the data provided by a relatively uniform participant set. While a rich understanding of the demographic

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2018 Association for Computing Machinery.

Manuscript submitted to ACM

and other personal characteristics of crowdworkers could play an important role in quantifying and subsequently addressing these biases, collecting such information comes with a number of drawbacks. For example, crowdworkers are repeatedly asked to provide the same information (*e.g.*, age, gender, income), increasing time on task, offering limited intellectual or creative stimulation, and providing little incentive to provide accurate details. Finally, from a task requester perspective, workers need to be financially compensated for their time spent on questions that are frequently repeated between tasks and are faced with negative consequences on data quality.

Recognising the importance of rich information on data labellers while simultaneously acknowledging the already existing financial and temporal pressure faced by crowdworkers, we present a conceptual browser plugin to assist crowdworkers in streamlining the data entry of repetitively requested personal information. This plugin can support the completion of demographic details, surveys on personality characteristics, and other data with a largely stable nature. We outline an infrastructure in which requesters can incorporate fast-form completion while gracefully degrading into an input form that can be used by crowdworkers not using such a browser plugin. Our envisioned system allows crowdworkers to focus on more exciting aspects of crowd work. To preserve crowdworker privacy and enable cross-platform compatibility, user data is to be stored on the worker’s browser. Given the extensive scale of contemporary crowdsourcing, this relatively small-scale improvement to everyday work can significantly impact the overall working conditions.

2 RELATED WORK

A multitude of HCI research has investigated different topics related to: fair pay for crowdworkers [3, 7, 15], empowering crowdworkers [10], or optimising task assignment to increase the quality of human intelligence task (HIT) results [9].

Hara et al. assessed the monetary compensation of crowdworkers based on a dataset of 3.8 million Amazon Mechanical Turk (MTurk) HITs from 2676 unique workers, showing that the average pay per hour is just \$3.13, with only 4.2% of the crowdworkers earning above the US minimum wage of \$7.25. These results are in stark contrast with existing calls for the fair compensation of crowdworkers by aligning to the minimum wage at the location of research [15]. Hara et al. highlight the existence of invisible work, which is work that does not result in earnings [7]. This includes, for example, the searching for HITs and rejected HITs. Rejected HITs are of particular interest, with rejected hits accounting for 12.8% of all submitted HITs and leading to an estimated monetary loss of 24% for crowdworkers [7]. Paying crowdworkers a minimum wage is not only ethically defensible, but may come with additional advantages. Litman et al. found that increasing the monetary reward for crowdworkers positively affected data quality as compared to a reduced compensation, with results holding true both in India and the US [13]. Further, Litman et al. found that the monetary gain was significantly more important than other factors in motivating participants to complete HITs.

A way to optimise task assignments in crowdworking is to take into consideration crowdworkers’ performance across different cognitive tasks (*e.g.*, Stroop test or N-back test). Hettiachchi et al. measured participants’ executive functions (*i.e.*, inhibition control, cognitive flexibility, and working memory), and assigned or recommended HITs in line with crowdworkers’ cognitive ability. Recommending based on cognitive performance led to a significant increase in the HIT results by 5%-20%. To remove crowdworkers’ ‘invisibility’, and highlight the ethical treatment of crowdworkers, Irani and Silberman identified common problems expressed by crowdworkers [10]. Among the most prominent issues identified was the feeling that work was rejected unfairly or arbitrarily, as well as the slow turnaround time of receiving payment for completed HITs. To address these issues, by making workers aware of requesters performance, Irani and Silberman developed the Turkopticon browser plugin. Using four metrics (communicativity, generosity, fairness, and promptness), Turkopticon attempts to make requesters performance more transparent, thereby empowering crowdworkers.

3 REPETITIVE DATA COLLECTION

As is evident from the existing literature, crowdworkers face numerous challenges while completing their tasks. One under-explored aspect is the repetitive nature of frequently requested data elements. This includes typical demographic information, such as age, gender, and ethnicity, but can also refer to more extensive survey instruments. For example, the widely used Big Five personality trait questionnaire by Goldberg (over 7000 citations on Google Scholar) consists of 50 questions [6], which assess neuroticism, extraversion, openness to experience, agreeableness, and conscientiousness. Due to its length, answering this questionnaire accurately requires a substantial amount of time. When asked to repeat this questionnaire multiple times, workers may experience frustration and data quality can be reduced due to an increase in inaccurate answers provided. Analysis of the Big Five traits over a four-year period revealed that these traits are relatively stable [2], making this a suitable instrument for workers to store their input for usage in later tasks.

Despite the challenges workers face in relation to repetitive data collection, the reasons for collecting this data are often critical. Crowdsourcing work is human work, and any biases that are present in individuals are therefore likely to reflect in the collected results – see, for example, prior work by Otterbacher, which highlights the biases introduced in the creation of metadata [14]. More diverse constellations of crowdworkers has furthermore shown to result in higher data quality as compared to less diverse crowdworker teams [16]. Therefore, a diverse set of crowdworkers is typically desirable, and the first step towards this goal is to obtain information on who are the crowdworkers completing the tasks.

The screenshot shows a web browser window with the address bar displaying 'crowdsourcing-task.com/task?worker_id=07'. The main content area is titled 'Crowdsourcing task example' and contains a form for personal information. The form has two input fields: 'Your age' with the value '42' and 'Your gender' with the value 'Female'. A blue 'Submit' button is located below the gender field. A modal popup window titled 'New information detected' is overlaid on the form. It contains the text 'Do you wish to save the following information to Fast-Form?' and lists the entered data: 'Age: 42' and 'Gender: Female'. At the bottom of the popup are two buttons: 'Submit without saving' and 'Submit and save'.

Fig. 1. Following a worker submitting a task, the Fast-Form popup window highlights to the worker that new personal data has been recognised and asks the user whether this can be stored.

4 FAST-FORM ENTRY: CONCEPTUAL OVERVIEW

Following our initial understanding of the challenges experienced in the repetitive entry of identical information in crowdworker tasks, we developed a concept around fast-form entry for crowdsourcing. Inspiration for this work comes in part from the autofill functionality that is present in most modern browsers, providing users with a near-instant ability to fill out their address or credit card information on online forms. Our concept of Fast-Form entry is built around three design requirements;

- **Minimise worker effort required.** Installation and setup of the system should be effortless, in which only a minimal amount of data is requested from the worker. This minimal amount of data (*e.g.*, worker gender and date of birth) can subsequently be used to introduce the worker to the system as they complete tasks. Other data contributions (*e.g.*, extensive personality surveys) are to be collected and stored upon the initial completion of these specific tasks as part of regular crowdworking tasks. Figure 1 indicates how Fast-Form engages with the user in this scenario.
- **Active consent to data input.** Rather than autofill input on tasks, we stipulate that workers actively consent to data being entered into an input field. This has the added benefit of the crowdworker being able to immediately validate that the entered data is correct, adding to the overall reliability of the crowdworker contributions.
- **Worker control over data storage and entry.** Given the personal and sometimes private nature of the data, it is critical that crowdworkers are in control of the data being shared [1]. We therefore follow a decentralised system design approach, in which data is stored on the client side and where the user has full control over who the data is shared with.

4.1 Envisioned infrastructure

Rather than expanding on the creation of detailed profiles of crowdworkers by crowdsourcing platforms¹, which eventually leads to lock-in and reduces workers' privacy, our envisioned infrastructure stores worker data on their

¹Prolific, for example, stores and provides information on country of birth, country of residence, and employment status among other demographic information.

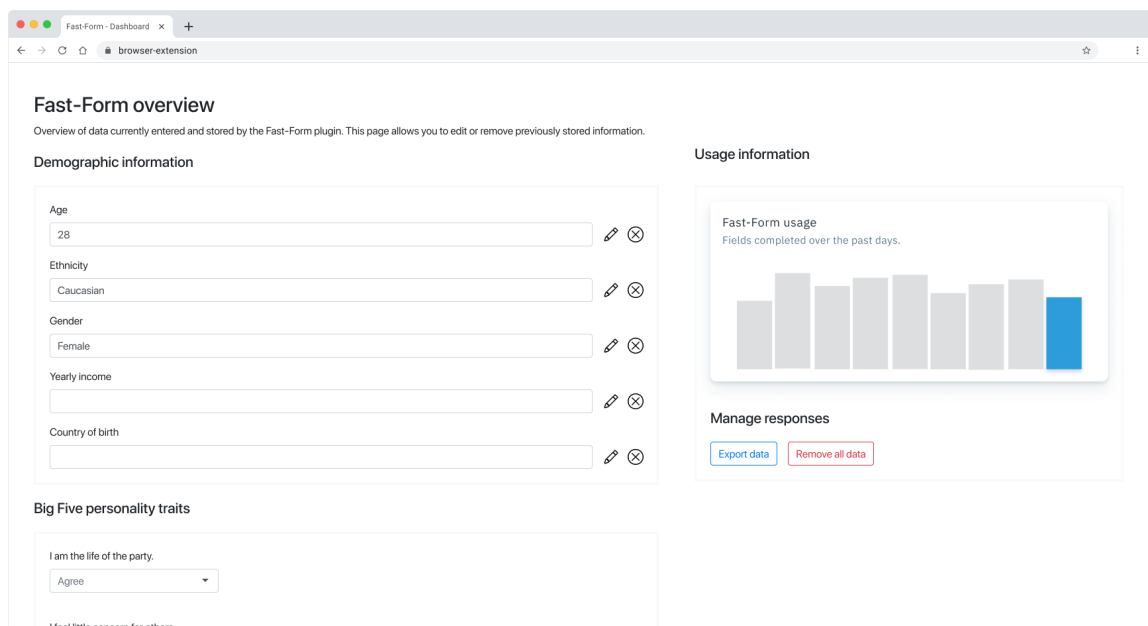


Fig. 2. Overview of worker data stored, providing the option to edit and remove previously recorded data. The plugin could also record frequency of use as an indicator of value to the worker.

personal machines (e.g. on chrome.storage [4]). The plugin configuration page should therefore allow crowdworkers to alter or remove stored personal data at any time, ensuring that the provided information is up-to-date and relevant.

To reduce worker effort required, we envision that requesters add a custom attribute to the input tags in their form input fields. The browser plugin automatically scans the task pages for elements with the custom attribute and subsequently alter the HTML's source code to provide the user interface elements that allow workers to provide fast-form entry answers. This method avoids interfering with the crowdsourcing platforms and keeps the logic separated from commonly used attributes such as *class* and *id*. For example, for the collection of a worker's age, an input field can be assigned a custom attribute with the name 'fastform' and a value of 'age' which triggers the plugin. By embedding this 'hook' in the task code itself, we ensure interoperability between different crowdworker platforms. See for example the following input field;

```
<input type="number" fastform="age">
```

An open challenge in our envisioned setup is the expiry of data. Various personal characteristics, such as those captured in the frequently used personality surveys, are likely to shift over long periods of time. To combat this, the plugin could, for example, remove data that has 'expired'. Such a consideration would require a careful and well-motivated balance between increased worker effort and expected gain in data quality. Continuously updating data based on future task completion may be able to alleviate some of these issues.

5 DISCUSSION AND CONCLUSION

In this paper, we presented a concept for the fast-form entry of frequently encountered survey data by crowdworkers. Crowdworkers face numerous challenges in their work [3, 7, 15], of which the repetitive nature of many identical survey tasks is under explored. In addition to reducing workers' frustrations, we argue that by reducing the barriers (e.g., time, costs) of repetitive data entry on stable personal characteristics, data requestors can collect richer background information on the individuals who provide them with crucial information. These insights can be crucial in identifying issues surrounding bias in datasets, which can form the basis of undesired discriminatory algorithms.

While our contribution is primarily conceptual in nature and does not consist of an evaluation, we anticipate that our work can serve as an interesting discussion point at the CHI REGROW workshop. Through this workshop, we hope to contribute to the growing realisation among the HCI community that we have not only the necessary knowledge and tools but also the obligation to ensure crowdworkers can perform their tasks in an enjoyable and satisfactory manner. Our concept highlights that such improvements may not only benefit the workers, but can ultimately also prove invaluable to task requestors.

REFERENCES

- [1] Andy Alorwu, Aku Visuri, Niels van Berkel, and Simo Hosio. 2022. (Re)Using Crowdsourced Health Data: Perceptions of Data Contributors. *IEEE Software* 39, 1 (2022), 36–42. <https://doi.org/10.1109/MS.2021.3117684>
- [2] Deborah A. Cobb-Clark and Stefanie Schurer. 2012. The stability of big-five personality traits. *Economics Letters* 115, 1 (2012), 11–15. <https://doi.org/10.1016/j.econlet.2011.11.015>
- [3] Greg d'Eon, Joslin Goh, Kate Larson, and Edith Law. 2019. Paying Crowd Workers for Collaborative Work. *Proc. ACM Hum.-Comput. Interact.* 3, CSCW, Article 125 (nov 2019), 24 pages. <https://doi.org/10.1145/3359227>
- [4] Google Developers Documentation. 2022. Chrome.storage. <https://developer.chrome.com/docs/extensions/reference/storage/>
- [5] Tarleton Gillespie. 2018. *Custodians of the Internet: Platforms, Content Moderation, and the Hidden Decisions That Shape Social Media*. Yale University Press, New Haven, Connecticut, USA. <https://doi.org/10.12987/9780300235029>
- [6] Lewis R. Goldberg. 1992. The development of markers for the Big-Five factor structure. *Psychological Assessment* 4, 1 (1992), 26–42. <https://doi.org/10.1037/1040-3590.4.1.26>

- [7] Kotaro Hara, Abigail Adams, Kristy Milland, Saiph Savage, Chris Callison-Burch, and Jeffrey P. Bigham. 2018. *A Data-Driven Analysis of Workers' Earnings on Amazon Mechanical Turk*. Association for Computing Machinery, New York, NY, USA, 1–14. <https://doi.org/10.1145/3173574.3174023>
- [8] Danula Hettiachchi and Jorge Goncalves. 2019. Towards Effective Crowd-Powered Online Content Moderation. In *Proceedings of the 31st Australian Conference on Human-Computer-Interaction (Fremantle, WA, Australia) (OZCHI'19)*. Association for Computing Machinery, New York, NY, USA, 342–346. <https://doi.org/10.1145/3369457.3369491>
- [9] Danula Hettiachchi, Niels van Berkel, Vassilis Kostakos, and Jorge Goncalves. 2020. CrowdCog: A Cognitive Skill Based System for Heterogeneous Task Assignment and Recommendation in Crowdsourcing. *Proc. ACM Hum.-Comput. Interact.* 4, CSCW2, Article 110 (oct 2020), 22 pages. <https://doi.org/10.1145/3415181>
- [10] Lilly C. Irani and M. Six Silberman. 2013. Turkopticon: Interrupting Worker Invisibility in Amazon Mechanical Turk. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems (Paris, France) (CHI '13)*. Association for Computing Machinery, New York, NY, USA, 611–620. <https://doi.org/10.1145/2470654.2470742>
- [11] Otto Kässi and Vili Lehdonvirta. 2018. Online labour index: Measuring the online gig economy for policy and research. *Technological Forecasting and Social Change* 137 (2018), 241–248. <https://doi.org/10.1016/j.techfore.2018.07.056>
- [12] Laura Lascau, Sandy J. J. Gould, Anna L. Cox, Elizaveta Karmannaya, and Duncan P. Brumby. 2019. Monotasking or Multitasking: Designing for Crowdworkers' Preferences. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems (Glasgow, Scotland Uk) (CHI '19)*. Association for Computing Machinery, New York, NY, USA, 1–14. <https://doi.org/10.1145/3290605.3300649>
- [13] Leib Litman, Jonathan Robinson, and Cheskie Rosenzweig. 2015. The relationship between motivation, monetary compensation, and data quality among US-and India-based workers on Mechanical Turk. *Behavior research methods* 47, 2 (2015), 519–528.
- [14] Jahna Otterbacher. 2015. Crowdsourcing Stereotypes: Linguistic Bias in Metadata Generated via GWAP. In *Proceedings of the 33rd Annual ACM Conference on Human Factors in Computing Systems (Seoul, Republic of Korea) (CHI '15)*. Association for Computing Machinery, New York, NY, USA, 1955–1964. <https://doi.org/10.1145/2702123.2702151>
- [15] M. S. Silberman, B. Tomlinson, R. LaPlante, J. Ross, L. Irani, and A. Zaldivar. 2018. Responsible Research with Crowds: Pay Crowdworkers at Least Minimum Wage. *Commun. ACM* 61, 3 (feb 2018), 39–41. <https://doi.org/10.1145/3180492>
- [16] Niels van Berkel, Jorge Goncalves, Danula Hettiachchi, Senuri Wijenayake, Ryan M. Kelly, and Vassilis Kostakos. 2019. Crowdsourcing Perceptions of Fair Predictors for Machine Learning: A Recidivism Case Study. *Proc. ACM Hum.-Comput. Interact.* 3, CSCW, Article 28 (nov 2019), 21 pages. <https://doi.org/10.1145/3359130>