

API vs. human: Comparing the performance of speech-to-text transcription using voice answers from a smartphone survey

Höhne and Lenzner

*DZHW, Leibniz University Hannover
GESIS – Leibniz Institute for the Social Sciences*

BigSurv23 Conference

Quito (Ecuador) – 26 to 29 October 2023



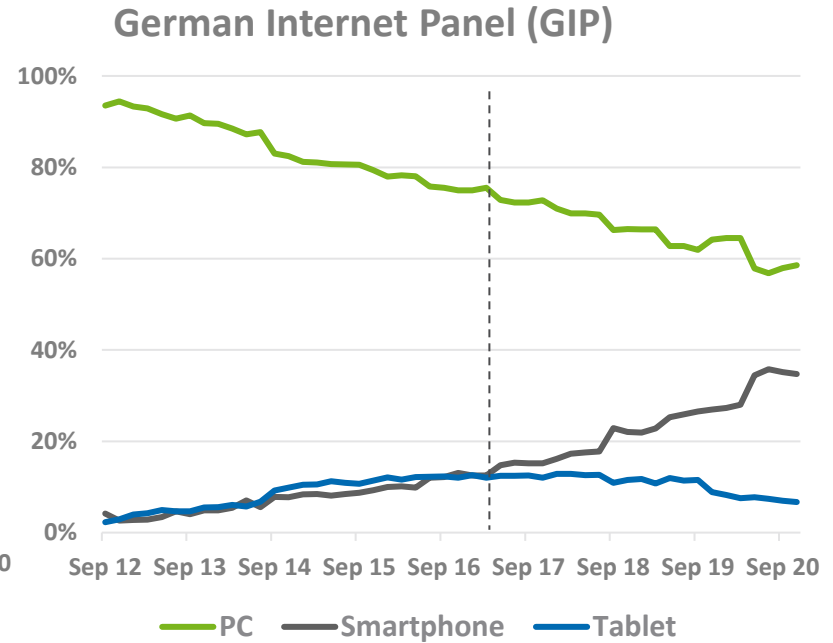
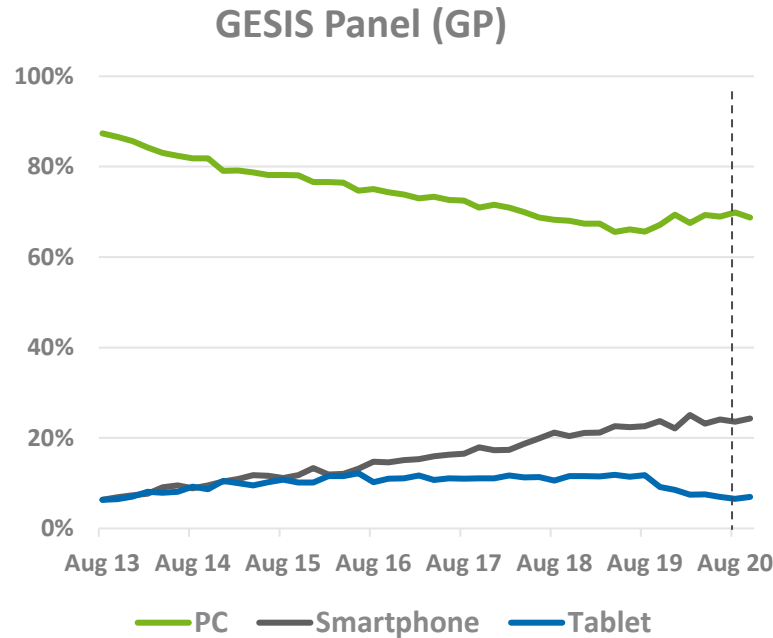
DZHW

Deutsches Zentrum für
Hochschul- und Wissenschaftsforschung

Web Surveys and Digital Innovations

- Increase of web-based surveys
 - *Academia: ANES, CRONOS, EVS, GESIS Panel, GIP, HRS, LISS Panel etc.*
 - *Public/private sector: Facebook, Google, UNESCO, World Bank etc.*
- Increase of mobile device use in web-based surveys
 - *Mobile optimized layouts as default* (Revilla et al. 2016)
- Emergence of digital intersections
 - *Ex ante data linkage (e.g., sensors)* (Elevelt et al. 2021; Höhne & Schlosser 2019)
 - *Ex post data linkage (e.g., trace data)* (Pasek et al. 2020; Stier et al. 2020)

Devices in Web Surveys



Country: Germany. Prob-based online panels (GP and GIP). Six waves per year. Vertical lines indicate the introduction of mobile-optimized layouts. Calculations: Gummer et al. (2023).

Smartphones and Voice Answers

- New communication channels because of smartphones
 - *Linking established methods with technological innovations*
- Voice answers to (open) questions
 - *Closeness to daily conversation* (Tourangeau et al. 2000)
 - *Rich information due to narrations* (Gavras & Höhne 2022; Gravras et al. 2022)
- Technological requirements of voice answers are met
 - *Even in web-based surveys with large N*
- Nonetheless, answer transcription is still required
 - *Human transcription is burdensome and time consuming*
 - *APIs may not be entirely ready*

Research Questions

- RQ1: What is the transcription quality of APIs?
- RQ2: What types of errors occur in API transcription?
- RQ3: How long does transcription by APIs and humans take?

Method: Study Design

forsa.omninet

Inwieweit stimmen Sie der folgenden Aussage zu oder nicht zu?

Ich fühle mich eher als Weltbürger und somit verbunden mit der Welt insgesamt und weniger als Bürger eines bestimmten Landes.

Stimme voll und ganz zu

Stimme zu

Weder noch

Stimme nicht zu

Stimme überhaupt nicht zu

Kann ich nicht sagen

[< Zurück](#) [Weiter >](#)

forsa. [Impressum](#) [Datenschutz](#)

forsa.omninet

Wie haben Sie den Begriff "Weltbürger" in der letzten Frage verstanden?

Halten Sie das Mikrofon-Symbol gedrückt, während Sie Ihre Antwort aufnehmen.

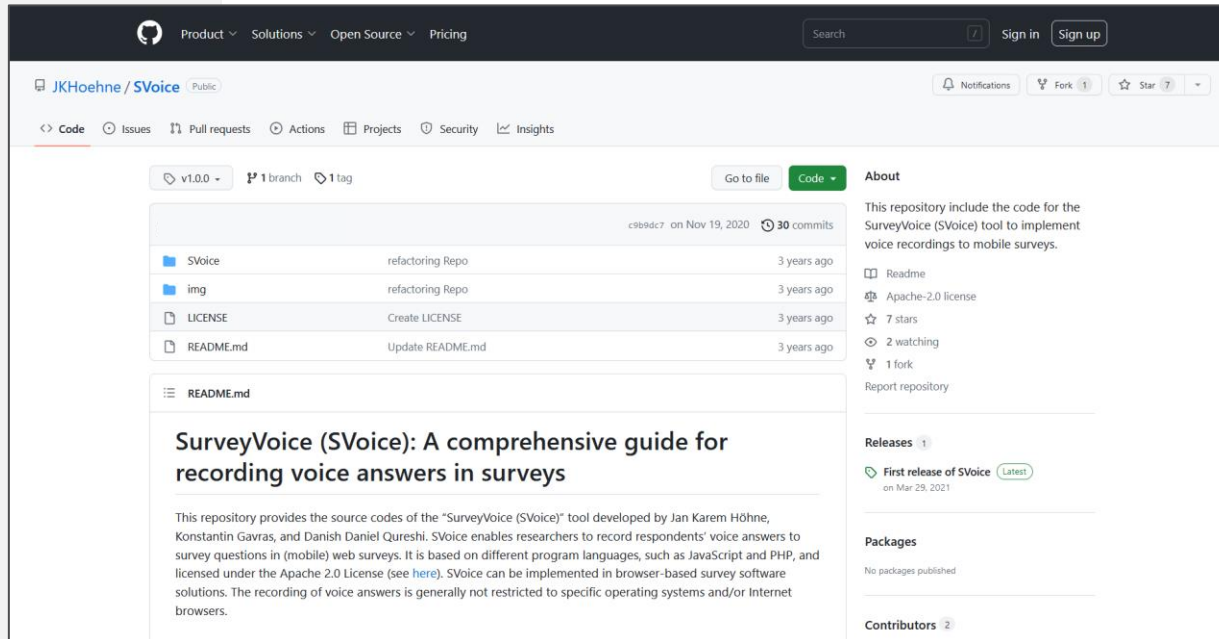


[< Zurück](#) [Weiter >](#)

forsa. [Impressum](#) [Datenschutz](#)

- Cross-quota sample
 - *Age, gender, and education*
 - *Forsa Omninet Panel (Nov 21)*
- 2 Questions + probes
 - *Relationship between citizens and state (ISSP 2013, 2014)*
- No recording time restrictions
 - *Overall, we have 609 voice answers for analysis.*
 - *These answers vary between 1 and 295 seconds.*

Method: Collecting Voice Data



The screenshot displays the GitHub interface for the repository `JKHoehne/SVoice`. The repository is public and has 7 stars and 1 fork. The main content area shows a file browser with the following files and folders:

- `SVoice` (refactoring Repo, 3 years ago)
- `img` (refactoring Repo, 3 years ago)
- `LICENSE` (Create LICENSE, 3 years ago)
- `README.md` (Update README.md, 3 years ago)

The `README.md` file is expanded, showing the following text:

SurveyVoice (SVoice): A comprehensive guide for recording voice answers in surveys

This repository provides the source codes of the "SurveyVoice (SVoice)" tool developed by Jan Karem Höhne, Konstantin Gavras, and Danish Daniel Qureshi. S.Voice enables researchers to record respondents' voice answers to survey questions in (mobile) web surveys. It is based on different program languages, such as JavaScript and PHP, and licensed under the Apache 2.0 License (see [here](#)). S.Voice can be implemented in browser-based survey software solutions. The recording of voice answers is generally not restricted to specific operating systems and/or Internet browsers.

The right sidebar contains the following sections:

- About:** This repository include the code for the SurveyVoice (S.Voice) tool to implement voice recordings to mobile surveys.
- Releases:** 1 release. **First release of S.Voice** (Latest) on Mar 29, 2021.
- Packages:** No packages published.
- Contributors:** 2 contributors.

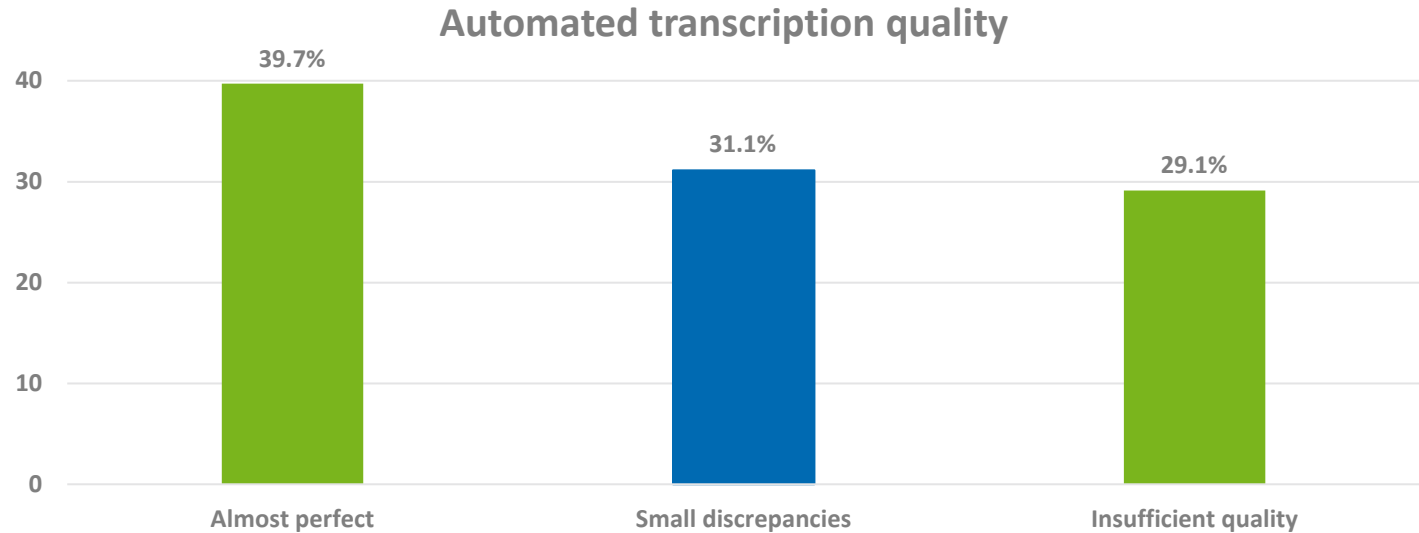
<https://github.com/JKHoehne/SVoice/tree/v1.0.0>

- SurveyVoice (S.Voice) tool (Höhne et al. 2021)
- Open-source
 - *Apache 2.0 License*
- JavaScript, CSS, HTML, and PHP
- Implementable in browser-based smartphone surveys

Method: Analytical Strategies

- We used the Google “Cloud Speech-to-Text” API
 - *Costs: \$0.024 per minute (without data logging – standard)*
- Transcription quality
 - *1 Almost perfect transcriptions, 2 small discrepancies (minor errors), and 3 insufficient quality (major errors)*
 - *Intercoder reliability: Agreement > 90% (Kappa > 0.84)*
- Transcription error types
 - *1 No mistakes, 2 misspellings, 3 word separation error, 4 word transcription error, 5 missing words, 6 incorrect grammatic, 7 words added by mistake, and 8 words replaced by numbers*
 - *Intercoder reliability: Agreement > 82% (Kappa > 0.78)*

Results: Research Question 1



We have less than 1% of poor-quality recordings that cannot be transcribed.

Results: Research Questions 1

Almost perfect

ja Weltbürger wäre dass ich mich überall zu Hause fühlen würde **das** oder so ähnlich

ich finde den Begriff **zivile** Ungehorsam doof und ich habe es auch dreimal gelesen um es einigermaßen zu verstehen **aber** ich will keine Beispiele nennen aber ich denke einfach dass wir anderer Meinung sind also die Menschen **da** anderer Meinung sind und schon deshalb einfach Ungehorsam aber der Begriff ziviler Ungehorsam der gehört hier nicht her widerstrebt mir finde ich auch nicht gut

Small discrepancies

ein Weltbürger ist quasi überall auf der Welt zu Hause vielleicht auch nirgendwo zu 100% zu Hause und ja **kann man** sich nicht an eine bestimmte Kultur oder Herkunft ist offen für **er** alle möglichen Kulturen und auch **ich freue mich** flexibel also möglichst viel gereist und hat auch ja einen großen Teil seines Lebens vielleicht im Ausland verbracht und dort gelebt

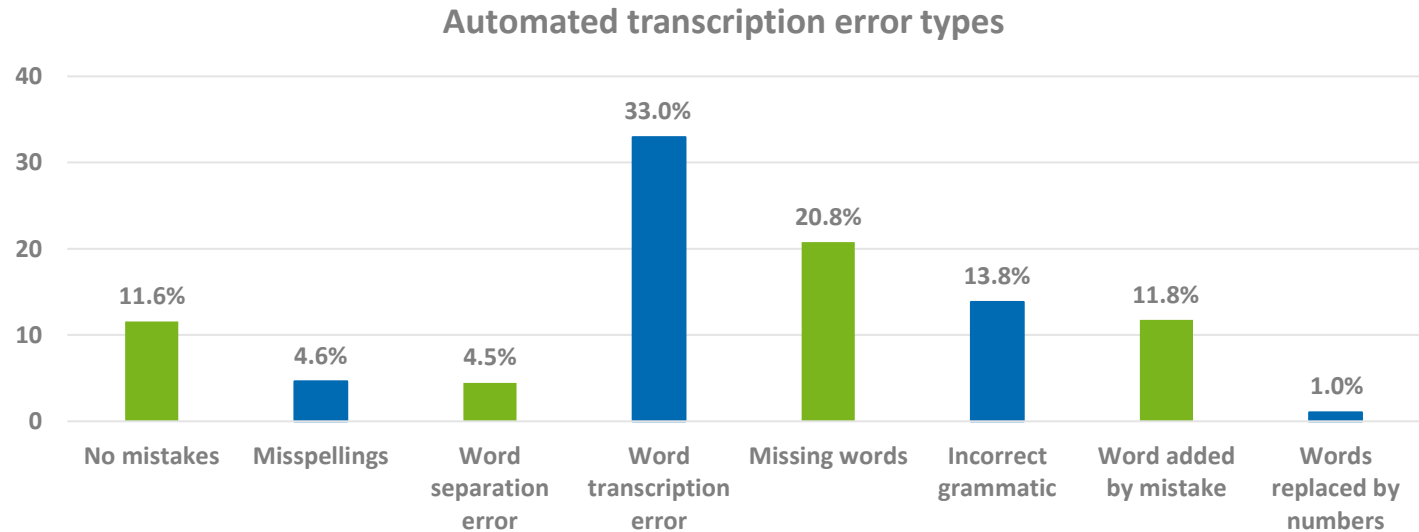
Beispiel ich bin zivil und gehorche nicht **in** Medien **und** Politik

Insufficient quality

es ist eine schöne **und mein** schöner **Garten Schüler glaube** oder sonst **du für** Weltbürger ist letztendlich wird die Mehrheit der Menschheit immer erst **zum 80.** somit auch auf ihr eigenes Land und **werde snap** macht der **geh** zu Grunde dafür gibt genügend Beispiele in der Geschichte

ja **heute** kleine Aufkleber **wie ging** das System also sich nicht an alles halten was die einem vorschreiben möchten weil es kann ja nicht sein dass das da das ist **zwar die** Meinungsfreiheit **los heute beginnt jetzt immer** Kasperl und so entscheiden **19 Zoll für dich** richtig ist aber naja **und das Badezimmer zwei Leute hatten da** ich glaube ich würde hier die **auch mal besprechen**

Results: Research Question 2



Results: Research Questions 3

Questions	File number	Coder 1 time (20%)	Coder 2 time (20%)	Mean time (20%)	Estimated time (100%)
1	63	125 min	95 min	110 min	550 min
2	60	120 min	100 min	110 min	550 min

Google “Cloud Speech-to-Text” API needs about 40 seconds for transcribing 1 min voice input. In total, the API needed about 153 minutes for all voice answers.

Discussion and Conclusion

- The API makes a great job in 40% of the cases
 - *In the other cases, there are (some) quality concerns*
 - *Combination of API and humans may be most promising*
- There are various error types that need to be considered
 - *Word transcription error and missing words are most prominent*
 - *Misspellings and word separation error are only minor threats*
- API transcription is 3.6 times faster than humans
 - *Its console can be easily operated through R*
 - *Open-source solutions are developing (e.g., Whisper)*
- Take home message
 - *APIs are not entirely ready for transcribing voice answers*
 - *Transcription effort needs to be considered in studies with voice answers*

Mobile Apps and Sensors in Surveys (MASS)



Call for Papers

5th Mobile Apps and Sensors in Surveys (MASS) Workshop

<https://massworkshop.org/>

Organizing committee

Jan Karem Höhne (DZHW, Leibniz University Hannover), Florian Keusch (University of Mannheim), Peter Lugtig (Utrecht University), and Bella Struminskaya (Utrecht University)

Date

March 6-7, 2024

Location

Washington, DC (USA)

This time, the Mobile Apps and Sensors in Surveys (MASS) workshop takes place jointly with the Current Innovations in Probability-based Household Internet Panel Research (CIPHER; see <https://dornsife.usc.edu/cesr/cipher-2024/>) conference. CIPHER will be held from March 7-8, and presenters at either event are invited to attend both events.

- MASS takes place in Washington DC (US)
 - From **March 6-7, 2024**
 - Together with CIPHER conference
 - Submission deadline is **December 1, 2023**



Many thanks for your attention!

www.jkhoehne.eu

@jkhoehne.bsky.social

@jkhoehne

Literature I

- Elevelt, A., Bernasco, W., Lugtig, P., Ruiter, S., & Toepoel, V. (2021). Where you at? Using GPS locations in an electronic time use diary study to derive functional locations. *Social Science Computer Review*, 39, 509–526.
- Gavras, K., & Höhne, J.K. (2022). Evaluating political parties: Criterion validity of open questions with requests for text and voice answers. *International Journal of Social Research Methodology*, 25, 135-141.
- Gavras, K., Höhne, J.K., Blom, A., & Schoen, H. (2022). Innovating the collection of open-ended answers: The linguistic and content characteristics of written and oral answers to political attitude questions. *Journal of the Royal Statistical Society (Series A)*, 185, 872-890.
- Gummer, T., Höhne, J.K., Rettig, T., Roßmann, J., & Kummerow, M. (2023). Is there a growing use of mobile devices in web surveys? Evidence from 128 web surveys in Germany. *Quality and Quantity*. DOI: 10.1007/s11135-022-01601-8
- Höhne, J.K., Gavras, K., & Qureshi, D.D. (2021). SurveyVoice (SVoice): A comprehensive guide for collecting voice answers in surveys. Zenodo. DOI: 10.5281/zenodo.4644590
- Höhne, J.K., & Schlosser, S. (2019). SurveyMotion: What can we learn from sensor data about respondents' completion and response behavior in mobile web surveys? *International Journal of Social Research Methodology*, 22, 379–391.
- Pasek, J., McClain, C.A., Newport, F., & Marken, S. (2020). Who's tweeting about the president? What big survey data can tell us about digital traces? *Social Science Computer Review*, 38, 633–650.

Literature II

- Revilla, M., Toninelli, D., Ochoa, C., & Loewe, G. (2016). Do online access panels need to adapt surveys for mobile devices? *Internet Research*, 26(5), 1209–1227.
- Stier, S., Breuer, J., Siegers, P., & Thorson, K. (2020). Integrating survey data and digital trace data: Key issues in developing an emerging field. *Social Science Computer Review*, 38, 503–516.
- Tourangeau, R., Rips, L.J., & Rasinski, K. (2000). *The psychology of survey response*. Cambridge, UK: Cambridge University Press.