

# Transcribing and coding voice answers obtained in web surveys: comparing three leading automatic speech recognition tools and human versus LLM-based coding

Melanie Revilla

*Research and Expertise Centre for Survey Methodology (RECSM),  
University Pompeu Fabra*

Carlos Ochoa

*Research and Expertise Centre for Survey Methodology (RECSM),  
University Pompeu Fabra*

Jan Karem Höhne

*German Centre for Higher Education Research and Science Studies (DZHW)  
Leibniz University Hannover*

Mick P. Couper

*University of Michigan*

## Abstract

With the rise of smartphone use in web surveys, voice or oral answers have become a promising methodology for collecting rich data. Voice answers not only facilitate broader and more detailed narratives but also include additional metadata, such as voice amplitude and pitch, to assess respondent engagement. Despite these advantages, challenges persist, including high item non-response rates, mixed respondent preferences for voice input, and labor-intensive, manual answer transcription and coding. This study addresses these last two challenges by evaluating two critical aspects of processing voice answers. First, it compares the transcription performance of three leading Automatic Speech Recognition (ASR) tools—Google Cloud Speech-to-Text API, OpenAI Whisper, and Vosk—using voice answers collected from an open-ended question on nursing home transparency that was administered in an opt-in online panel in Spain. Second, it evaluates the efficiency and quality of coding these transcriptions using human coders and GPT-4o, a Large Language Model (LLM) developed by OpenAI. We found that each of the ASR tools has distinct merits and limits. Google sometimes fails to provide transcriptions, Whisper produces hallucinations (false transcriptions), and Vosk has clarity issues and high rates of incorrect words. Human and LLM-based coding also differ significantly. Thus, we recommend using several ASR tools for voice answer transcription and

---

This document is a preprint and thus it may differ from the final version.

Data availability: The anonymized datasets and R scripts used for their analysis are accessible in Open Science Framework (OSF): [https://osf.io/8rjb6/?view\\_only=5e3156ec585641ccb22b0b712c0af0cf](https://osf.io/8rjb6/?view_only=5e3156ec585641ccb22b0b712c0af0cf). All the supplementary online materials (SOM) are also available in the same folder.

Funding: This work was supported by the European Research Council (ERC) under the European Union's Horizon 2020 research and innovation program [grant number 849165].

Acknowledgements: We are very grateful to Maria Paula Acuña Pardo for her hard work on the coding, to Ixchel Perez Duran for her work in designing the questionnaire, and to Lucía Fernández Melero, Maiki Estevez Cano and Joshua Claassen, for their assistance at different stages of this manuscript preparation.

implementing human as well as LLM-based coding, as the latter offers additional information at minimal added cost.

*Keywords: Automatic Speech Recognition (ASR), Google's Cloud Speech-to-Text API, GPT-4o, Large Language Model (LLM), Voice answer transcription, Vosk, OpenAI Whisper*

## **1. Introduction**

Web surveys offer notable benefits for both respondents (flexibility in terms of time and location) and researchers: especially, timeliness, cost-efficiency (Callegaro et al. 2015), and technological adaptability (Conrad et al. 2021; Struminskaya et al. 2020), enhanced by growing smartphone participation rates (Gummer et al. 2019, 2023; Peterson et al. 2017; Revilla et al. 2016). In particular, built-in microphones facilitate the administration of voice or oral answers to open-ended questions (Revilla 2022; Schober et al. 2015).

Requesting voice answers in web surveys offers potential benefits, again for both researchers and respondents (Revilla 2022). From a researcher's perspective, answers gathered through voice recording in web surveys enable the collection of comprehensive information by triggering open narrations (Gavras & Höhne 2022), allowing respondents to articulate their thoughts more freely. Compared to written answers, oral answers tend to contain more words and characters (Gavras et al. 2022; Höhne & Claassen 2024; Revilla et al. 2020), while also taking less time (Revilla et al. 2020). Moreover, oral answers encompass a broader range of topics than written ones (Gavras et al. 2022) and are associated with higher levels of validity (Gavras & Höhne 2022). The meta-data included in oral answers, such as voice amplitudes and pitches, can be utilized to gauge respondents' interest levels during survey completion (Höhne et al. 2024), introducing an additional dimension to the analysis of answer behavior and data quality. From a respondent perspective, answering through voice recordings might be faster and less burdensome than typing in a text box, especially on smartphones. Additionally, this method may feel more natural and enjoyable, as many people regularly use voice functions, including those of voice assistants and instant messengers, in their daily lives (Deloitte 2018; Revilla et al. 2018).

Despite these advantages, previous research reports item non-response rates between 25% and 60% (Gavras et al. 2022; Revilla et al. 2020; Revilla & Couper 2021). Moreover, even though most respondents find it easy to answer through voice recording, only 39% reported liking it (Revilla & Couper 2024). Another hurdle linked to oral answers is the need to transcribe them into text for substantive analysis. In theory, Automatic Speech Recognition (ASR) systems can process voice input directly and automatically code this input to trigger specific actions, such as in Customer Relationship Management (CRM) systems. However, these systems usually have limited functionality and can only directly interpret a narrow range of inputs (e.g., brief commands or simple "yes" or "no" answers). Thus, at this time, they lack the ability to handle complex, narrative answers. Consequently, such systems typically fail to extract the detailed information researchers seek from narrative open-ended answers. Coding such answers often requires careful consideration of multiple aspects and may involve repeated review. Therefore, a two-step process – first transcribing the answers and then coding them – is usually necessary.

This requirement introduces an additional stage to data processing that must be accounted for, as transcriptions typically entail substantial time and personnel effort. Transcribing audio files typically takes three to eight times longer than the original voice input (McMullin 2023). ASR tools present a potential solution to bypass effortful manual transcription through automatic transcription. Despite claims of their effectiveness across various languages, there are limited empirical demonstrations testing ASR tools, particularly within the context of web surveys (see Section 2).

Once transcribed into text, the oral answers can be coded similarly to conventional written answers. While human coding has often been employed (Höhne & Claassen 2024; Höhne et al. 2025; Lenzner et al. 2024; Revilla et al. 2020, Revilla & Couper 2021), it is labor- and time-intensive and inconsistencies can arise as different coders may assign different codes to the same answers. To address these challenges, some studies focusing on open-ended narrative questions have opted for automatic coding approaches utilizing Natural Language Processing (NLP) and machine learning, such as Structural Topic Modeling (STM; Roberts et al. 2014) and Bidirectional Encoder Representations from Transformers (BERT; Landesvatter & Bauer 2024). However, these methods have their own limitations. For example, STM does not consider the word order and grammatical structure of responses, which may lead to inaccuracies or a lack of interpretive depth (Barde & Bainwad 2017). BERT models require a sufficient number of open-ended answers for model fine-tuning (Gweon & Schonlau 2024).

Recent advances in Large Language Models (LLMs) introduce new possibilities for survey researchers and practitioners for coding open-ended answers. In particular, GPT-4o, the most recent LLM of OpenAI when this research was conducted, can generate coherent text based on user input. This allows sophisticated and efficient answer coding, potentially reducing coding time, while maintaining interpretative richness (OpenAI et al. 2024). Some studies already compared GPT to human-based coding (see Section 2.3). However, empirical evidence remains limited. To our knowledge, a comparison of human and GPT-based coding of automatically transcribed voice answers has not been conducted yet. Thus, this study has two main goals:

1) Providing new empirical evidence about the performance of three leading ASR tools to transcribe voice answers to an open-ended question on nursing homes transparency that were collected through a web survey administered in the Netquest online panel<sup>1</sup> in Spain in 2024: Google's Cloud Speech-to-Text API<sup>2</sup>, OpenAI Whisper<sup>3</sup>, and Vosk<sup>4</sup>. Google and Whisper were chosen because these are the two ASR tools that have been compared in previous research on voice answer transcription (Höhne et al. 2025; Meitinger et al. 2024). Vosk was included because it has been used in previous studies to transcribe voice answers for substantive analysis (Revilla and Couper 2023). Additionally, all three tools offer several advantages (see Section 2.1), including low or no costs and Spanish language capability.

2) Comparing different ways to code these automatically transcribed voice answers: human and LLM-based coding (GPT-4o model), using two different parameter settings (see Section 3).

---

<sup>1</sup> [www.netquest.com](http://www.netquest.com)

<sup>2</sup> <https://cloud.google.com/speech-to-text>

<sup>3</sup> <https://openai.com/index/whisper/>

<sup>4</sup> <https://github.com/alphacep/vosk-api/>

## 2. Background

### 2.1 Performance of Google, Whisper, and Vosk

When evaluating ASR tools, key factors, such as accuracy under various conditions (e.g., noise levels and accents), typically measured by Word Error Rate (WER) and processing speed, are essential. These criteria help to assess how well different ASR tools handle diverse audio inputs and environmental conditions.

The Google Cloud Speech-to-Text API is a commercial tool powered by its Universal Speech Model, which relies on a family of advanced speech models with two billion parameters, trained on 12 million hours of speech data and 28 billion sentences in over 300 languages ([Gladia](#)<sup>5</sup>). This extensive training enables it to excel in handling diverse accents and languages, currently supporting more than 125 languages ([Cloud Compiled](#)<sup>6</sup>). Furthermore, Google’s API performs well even when background noise is present, and is highly customizable through features like model adaptation, which allow it to recognize domain-specific terminology, enhancing its flexibility and accuracy for specialized applications ([Gladia](#)).

Whisper has been trained on 680,000 hours of multilingual and multitask data from online sources. Although the model was initially trained on 98 languages, only 50 languages (with a WER lower than 50%) are currently available ([Slator](#)<sup>7</sup>). Whisper can operate both locally on devices without internet access and online via its API. It has gained recognition for its accuracy, especially for difficult audio with background noise or multiple speakers and languages. Research found that Whisper performs well in terms of WER (Radford et al, 2023), outperforming both the Google API (Chen et al., 2024) and Vosk (Trabelsi et al., 2024), particularly in challenging scenarios. It also offers quicker processing time compared to the Google API when using Whisper’s model “Small” through the API (Chen et al., 2024). Using Whisper locally is free of charge, and its API is usually cost-effective for “smaller projects” ([Kenility](#)<sup>8</sup>). However, Whisper may not perform as well as other specialized models for cleaner, simpler datasets and is prone to “hallucinations,” as the model sometimes inserts extraneous words or phrases not present in the audio ([Gladia](#); [Slator](#)).

Vosk is an open-source ASR tool that leverages deep learning models along with optimized feature extraction techniques to transcribe audio into text. Unlike many cloud-based ASR tools, Vosk is designed to operate locally offline ([Medium](#)<sup>9</sup>). This makes it ideal for privacy-focused applications or where connectivity is unreliable. Thus, Vosk is widely used in offline scenarios and is praised for its efficiency and flexibility. It offers a good balance between performance and cost, although it may not reach the high accuracy levels of other tools for more complex or noisy audios. However, compared to Whisper, Vosk has been found to require fewer manual adjustments to ensure transcription accuracy ([Toolify](#)<sup>10</sup>).

### 2.2 Testing ASR tools using survey answers provided through voice

When focusing on answers to open narrative survey questions collected through voice recording, very few empirical studies exist. However, voice answers by survey respondents may

---

<sup>5</sup> <https://www.gladia.io/blog/openai-whisper-vs-google-speech-to-text-vs-amazon-transcribe>

<sup>6</sup> <https://cloudcompiled.com/2020/07/28/transcription-api-comparison/>

<sup>7</sup> <https://slator.com/resources/is-whisper-the-best-speech-to-text-software/>

<sup>8</sup> <https://www.kenility.com/blog/technology/rise-ai-transcription-whisper-vs-google-speech-text>

<sup>9</sup> <https://fahizkp.medium.com/vosk-a-comprehensive-guide-to-open-source-speech-recognition-3e634fc8d713>

<sup>10</sup> <https://www.toolify.ai/ai-news/enhanced-audiototext-comparison-vosk-vs-whisper-in-subtitle-edit-55557>

differ from the audio material usually used to compare the performance of ASR tools (see section 2.1). Several parameters that may impact performance are not under the researcher's control, including volume, speed of speech, accent, tone, or lexical structure. In particular, voice answers can be affected by background noise – respondents can answer whenever and wherever they want (Mavletova, 2013) – potentially lowering transcription accuracy (Pentland et al., 2023). Additionally, voice answers from smartphones are comparatively short – sometimes lasting only a few seconds – but ASR performance improves with the speech input length (Proksch et al., 2019).

Meitinger et al. (2024) explored the transcription accuracy of oral answers from the Longitudinal Internet Studies for the Social Sciences (LISS) panel in the Netherlands. Employing the [Questfox](#) tool<sup>11</sup>, that uses the Google API (transcription took place in 2020), they observed that background noise and the presence of third parties compromised transcription accuracy. Respondent characteristics, such as age and education, were not associated with transcription accuracy. Höhne et al. (2025) investigated the performance of Google API and Whisper using oral answers from a German non-probability online panel (transcription took place in 2024). In contrast to Chen et al.'s (2024) finding that Whisper is faster than Google, they report that Google processed and returned transcripts faster than Whisper (operated locally). However, this speed comes at the cost of more errors. The Google API produced around 20% transcriptions of insufficient quality with major errors, versus around 5% for Whisper.

The studies conducted by Meitinger et al. (2024) and Höhne et al. (2025) provide key empirical evidence. However, they considered a limited set of languages and ASR tools.

### ***2.3 Comparing human versus Generative Pre-trained Transformer (GPT) coding***

Recent studies have compared human and automated coding through OpenAI GPT-3.5 and GPT-4o models, highlighting their advantages and limitations. Automated coding with these models offers notable benefits in terms of accuracy, efficiency, and replicability (Theelen et al., 2024; Arlinghaus et al., 2024; Liu & Sun, 2023). Research has demonstrated the ability of these models to code and identify themes and patterns with high precision, often surpassing human coders in agreement rates, uncovering nuanced insights, and reducing bias while maintaining neutrality and consistency (Fuller et al., 2024; Liu & Sun, 2023). Furthermore, using these models significantly reduces the time and resources required for coding (Arlinghaus et al., 2024; Fuller et al., 2024).

However, OpenAI GPT also presents qualitative coding limitations, especially with techniques like axial coding (Saldaña, 2015), where tailored and refined prompts (i.e., written instructions for the model to guide its responses) are required to enhance GPT's task understanding (Theelen et al., 2024). OpenAI GPT models often require extensive context to produce meaningful codes that align with underlying theories (Fuller et al., 2024). Additionally, they may overgeneralize themes or overlook implicit nuances and emotions that human coders usually recognize (Liu & Sun, 2023).

Overall, previous research suggests that OpenAI GPT coding might be a promising solution, but important limitations remain. Additionally, while human coding can largely vary across coders, leading to low interrater reliability (IRR), LLM-based coding can also exhibit

---

<sup>11</sup> <https://questfox.online/en/questmanagement>

variability. Performance depends first on the specific LLM used. In particular, comparisons between GPT-3.5 and GPT-4o indicate that GPT-4o produces better explanations and higher agreement (Arlinghaus et al., 2024; Lee et al., 2024). Even when using the same LLM, the performance can vary depending on the exact coding task, the language of the text to be coded, and the formulation of the prompts. Moreover, GPT-based coding relies on several parameters, especially the “temperature,” which influences the model’s “creativity” or “focus.” A lower temperature produces more deterministic and focused outputs, increasing the likelihood of generating consistent coding when using the same prompts and data. In contrast, a higher temperature results in more creative and diverse outputs but reduces reproducibility (Marion, 2024).

Furthermore, new challenges may arise when applying GPT-based coding to ASR-based transcriptions. Indeed, as discussed earlier, ASR tools can introduce errors, such as misinterpreted words, grammatical issues, or even hallucinated content. These issues can exacerbate the limitations of GPT coding: inaccurate or poorly organized transcriptions may hinder GPT’s ability to produce accurate and meaningful coding outputs. Fuller et al. (2024) highlight the importance of data cleaning and crafting effective prompts. Untidy transcriptions from ASR tools, such as transcriptions that include spelling or punctuation errors and inconsistent capitalization, can lead to coding errors and inaccuracies, undermining the accuracy of coding results.

### 3. Research questions and contribution

This study’s first research question investigates the effectiveness of three leading ASR tools (Google Cloud Speech-to-Text, OpenAI Whisper, and Vosk) in transcribing oral answers (in Spanish) from web surveys:

***RQ1:** How do the ASR tools perform across various dimensions?*

By addressing **RQ1**, we extend the work by Meitinger et al. (2024) and Höhne et al. (2025). First, we consider another language. While Meitinger et al. (2024) considered Dutch and Höhne et al. (2025) considered German, we consider Spanish, a widely spoken language for which many large language datasets exist, which eases creation of ASR tools. Second, we expand the number of ASR tools under investigation, by not only exploring Google’s Cloud Speech-to-Text API and OpenAI Whisper, but also Vosk. Third, we use more recent versions of the tools. Specifically, in the study by Meitinger et al. (2024), data collection and oral answer transcription occurred in 2020. Given the rapid evolution of ASR tools, this raises questions about the ongoing validity of the results. Importantly, this is an exploratory evaluation, so we do not have specific expectations about the relative performance of the ASR tools.

We consider different aspects of performance of the ASR tools: whether a transcription is obtained, the number of characters, words, and sentences in the transcriptions, their clarity (i.e., how understandable and readable are transcriptions), presence of different kinds of problems, and validity of the answers (i.e., they align with the question and provide substantive information; see Section 4.4). Differences in the words used and meaning across pairs of ASR tools are also investigated. A key limitation of this study is that we do not have direct access to the original audio files of respondents, a restriction implemented to minimize data protection

risks. Thus, we cannot determine the “true” values (what respondents actually said), limiting performance evaluation to self-evident aspects. For example, incomplete sentences can be identified if they only include a subject without a verb, but a missing adjective (e.g., “very”) cannot be detected, as the sentence remains functional.

The second research question investigates different coding procedures:

***RQ2.** How similar or different are the codes of transcribed responses generated by a human and the OpenAI GPT-4o model?*

By addressing **RQ2**, we contribute to the limited body of research comparing human and LLM-based coding. We particularly focus on the respective performance in coding information from transcriptions of oral answers to narrative questions in web surveys. Since the results of GPT-4o depend on the parameter configuration, in particular the temperature<sup>12</sup>, we compare two GPT-4o outputs: one using the default temperature (0.7), and another setting the temperature to 0.

#### **4. Method and data**

This study uses a subset of the data collected in the framework of a pre-registered study investigating:

1) Whether the provision of extra incentives – given beyond the baseline incentive all respondents receive for their participation – to those answering experimental questions through voice increases the proportion of voice answers across groups of respondents varying in their likelihood of using voice.

2) The association between these additional incentives and the quality of answers.

For information about the full study design, we refer to Höhne, Revilla and Couper (2024). This section focuses on the aspects relevant to the current study.

##### **4.1 Questionnaire**

The questionnaire included over 80 questions, administered via a web survey optimized for mobile devices but also accessible on PCs. Due to routing, no respondent answered all questions. The full questionnaire and its English translation are available in SOM1. Respondents could skip questions, except those controlling quotas or tailoring subsequent questions.

The survey primarily focused on citizens’ perceptions of nursing homes in Spain but also included questions on political opinions, respondent characteristics, and web survey completion, among others.

This paper focuses on one open-ended narrative question<sup>13</sup> (see question “WHYTRANSP\_EXP” in SOM1) in which respondents were asked to explain why they

---

<sup>12</sup> Other settings could be adjusted, primarily: max\_tokens (Limit on response length), top\_p (Nucleus sampling, also known as “cumulative probability”), frequency\_penalty (Penalty for token frequency), presence\_penalty (Penalty for token presence), stop (Stop sequences), and logit\_bias (Adjustment of probabilities for specific tokens).

<sup>13</sup> Although the survey included a second open-ended narrative question with a request for voice answer, we focus on the data from WHYTRANPS\_EXP, because we do not expect differences in transcription performance between the two questions, as they share the same structure and topic. Additionally, Höhne et al. (2025) did not find differences in performance across their two questions.

selected a given answer in a prior closed question on the amount of information they think nursing homes in Spain provide.

For this question, a push-to-voice design was employed, where participants were initially asked to answer through voice recording. In a follow-up, respondents skipping the question were offered two options: record their answer or type it in a text-box. Since no differences are expected in the transcription performance, to reduce coding time and effort, we do not analyze the voice answers from the follow-up.

#### **4.2 Data collection**

Data was collected in the [Netquest](#) opt-in online panel in Spain between February 29, 2024, and March 22, 2024. To record respondents' oral answers, the *WebdataVoice* tool (Revilla et al. 2022), that works across devices (PCs, tablets, and smartphones) and mobile operating systems (Android and iOS), was used. Respondents were able to listen to their recordings before submitting them, and delete and re-record if needed. To minimize data disclosure risks, Netquest immediately transcribed the oral answers into text using the three ASR tools. These transcriptions were then forwarded to the projects' Ethics Advisor for manual review. In very few cases where unsolicited personal information was present, the advisor removed this information before sharing the final dataset with the research team.

We used quotas for gender and age (crossed) and education, to match the adult online population in Spain (under 75) according to the National Statistics Institute (see SOM1). Of the 11,076 panelists invited to the survey, 3,237 started it but 286 abandoned the survey before getting to the question of interest in this study and another 689 were excluded (e.g., for not giving their explicit consent to participate or for exceeding the quotas). Overall, 2,262 panelists got to the open-ended question under investigation. Of those, 1,403 panelists did not have any transcriptions, indicating they either initially skipped the question or encountered issues with their voice files. This leaves 859 panelists for our statistical analyses (those with at least one transcription).

The average age of these 859 panelists is 48 years, 51% of them are female, and 36% have a higher education degree. On average, they have been in the Netquest panel for 6.7 years (median = 6.4) and have completed 195 surveys (median = 170). About 21% completed the survey with a PC, 2% with a tablet, and 78% with a smartphone.

#### **4.3 Transcriptions**

The transcriptions of the audio files were done with the three ASR tools in September 2024<sup>14</sup>. For each tool, specific decisions were required regarding the models and parameters used. Detailed information about the configuration is provided in SOM2. Variations in the configuration can lead to different transcription outcomes, especially for the Google API. For this tool, the wide range of configuration options makes it more challenging to identify a set of parameters performing consistently well across audio files. Netquest observed that certain audio files failed to generate transcriptions with some settings but succeeded when adjusting the settings. Conversely, files that produced transcriptions with the initial settings sometimes failed

---

<sup>14</sup> Initial transcriptions were conducted immediately after data collection. However, not all audio files were included. Upon detecting this issue, Netquest implemented again all transcriptions: we use these new transcriptions.



after changes were made. Thus, Netquest tested multiple configurations for the Google API on a small subset of audio files and selected the settings that delivered the best performance metrics. In contrast, Netquest used the default settings for both Whisper and Vosk, as they performed reasonably well with these configurations. Since this part of the work was outsourced, our ability to control the process was limited.

#### 4.4 Coding of the transcriptions

We extracted various relevant information from the transcriptions, as presented in Table 1. The coding was organized in two blocks: Block 1 focuses on aspects coded for each ASR tool individually. Block 2, in contrast, addresses aspects that directly compare pairs of ASR tools (“Google-Whisper,” “Whisper-Vosk,” and “Vosk-Google”).

Table 1. Aspects coded

Block	Aspect coded	Coding method	What was coded
1	<i>Transcription provided</i>	R script	Binary variable. 1 indicates a transcription is provided. Any form of answer is considered, even nonsensical ones. We only consider respondents with a transcription for at least one of the three ASR tools, because we are interested in the relative performance of the three tools. The remaining aspects are coded for cases where an answer was observed with a given ASR tool.
1	<i>Answer length</i>	R script	Measured using three metrics: number of 1) characters, 2) words, and 3) sentences.
1	<i>Clarity</i>	Human & GPT-4o	Evaluates how understandable and readable the transcribed text is. Three levels are considered: 1) content is largely unclear (“not clear at all”), 2) some parts are unclear, but the overall transcription is usable (“clear”), and 3) content is very clear (“very clear”).
1	<i>Presence of different types of problems</i>	Human & GPT-4o	Three binary variables are used to code problems: 1) missing words or incomplete sentences, 2) words added by mistake or part of the answer is repeated, and 3) misspelling, grammatical errors, or wrong words (i.e., that do not make sense in the context of the text).
1	<i>No problem</i>	R script	Binary variable with value 1 if the three previous measures (missing, added, and wrong words) are 0, and value 0 otherwise.
1	<i>Valid answers</i>	Human & GPT-4o	Following Revilla and Couper (2023), we evaluate the validity of each transcription. Nonvalid answers include nonsense, answers not in line with the question topic, and non-substantive answers (e.g., “don’t know” or “no opinion”).

Table 1. *Continued*

Block	Aspect coded	Coding method	What was coded
2	<i>Number of different words</i>	R script	For each respondent, we count the number of words that differ between each pair of ASR tools (i.e., how many words appear in only one of the two transcriptions?). For example, if Vosk transcribes “I have a car” and Whisper transcribes “I have a bike,” the difference is two, as each transcription contains a unique word (car and bike, respectively).
2	<i>Percentage of different words</i>	R script	To contextualize these differences in relation to the overall length of each answer, we also compute, for each respondent and pair of ASR tools, the percentage of different words, by dividing the number of different words by the total number of words in both transcriptions combined and multiplying by 100. For the example on car and bike, it would be $2/8*100=25\%$ .
2	<i>Similarity of meaning</i>	Human & GPT-4o	Assesses how closely the meanings of each pair of transcriptions align. Three levels are considered: 1) the meaning of the transcriptions differs a lot for some parts (“not similar at all”), 2) the meaning is not exactly the same, but it is quite similar (“partly similar”), and 3) the meaning is identical (“very similar”).

Objective aspects (e.g., number of words) were extracted using R version 4.3.1 (R Core Team, 2023; script available in OSF). For subjective aspects, we use (and compare) both human coding and GPT-4o coding, fixing the temperature to 0 and using the default setting (0.7). Human coding was done in the last trimester of 2024 by a native Spanish speaker, following detailed guidelines (see SOM3). GPT-based coding took place in January and February 2025. Before conducting the full GPT-based coding, we first tested different approaches on a small subset of transcriptions, focusing primarily on prompt formulation. We began with prompts that closely matched the human coding guidelines, using the same examples. As we identified issues, we refined the prompts accordingly, especially to improve alignment between GPT’s decisions and those of the human coder in cases where a specific aspect was not explicitly addressed in the human coding guidelines. For instance, when assessing clarity, we added the following instruction, because GPT initially tended to classify transcriptions as unclear when only minor punctuation errors were present: “Do not lower your rating if the transcription has only minor punctuation issues”.

We also experimented on a subset of cases with different ways of sending the data to GTP-4o. First, we sent it all at once, expecting this approach to be simpler, faster, and more cost-effective, as pricing depends on the length of the texts sent to GPT. Thus, repeating the same prompt for each transcription increases costs. However, this method led to a reduced

proportion of coded transcriptions. For instance, if we sent 30 transcriptions at once, GPT failed to provide a code for some of them. Consequently, we opted to process the data incrementally, sending three transcriptions at a time and repeating the prompt as needed until all transcriptions were coded. Since we still had missing codes, we ultimately decided to send the requests one at a time. The final codes were generated by GPT-4o in 228 minutes with a temperature of 0 and in 219 minutes using the default temperature. The Python code (including the exact prompts) used for coding with GPT-4o can be found in SOM4.

#### 4.5 Analyses

The analyses were performed using R 4.3.1 (R Core Team, 2023; script available in OSF).

To answer **RQ1**, we evaluate each ASR tool for the different aspects coded in Block 1. We implement descriptive analyses of aspects that help to assess the quality of each transcription. For numeric variables, we report means, while for categorical and binary variables, we report proportions, expressed as percentages. With the exception of the first indicator (Transcription provided), we report the results for all participants for whom a transcription was obtained (Transcription provided = 1). However, given the substantial differences in the number of cases with a transcription across the ASR tools, we also conduct additional analyses focusing exclusively on cases where all three ASR tools produced a transcription. We then compare these tools: 1) By computing the differences between means or proportions for all pairs of ASR tools (“Google-Whisper,” “Whisper-Vosk,” and “Vosk-Google”) for each of the aspects in Block 1, and test whether these differences reach significance. 2) By analyzing the aspects in Block 2 that directly compare the ASR tools. For the number and percentages of different words, we report the average over all respondents with a transcription for each pair of ASR tools.<sup>15</sup> In SOM5, we report the results when focusing only on those with a transcription for all three tools. For meaning similarity, we report the proportions of transcription pairs with either partial or full meaning overlap. We conduct t-tests when comparing means and McNemar tests when comparing proportions. We report significant differences at the 5% level.

To address **RQ2**, we compare human and GPT-4o coding, when using the default temperature and temperature 0. Lowering the temperature is expected to reduce variation in coding but may also limit the model’s creativity, potentially leading to lower performance, particularly in a complex coding task, such as the one examined in this study. The comparison is first conducted by replicating the analyses from **RQ1**, excluding the variables created directly in R, and testing whether the results change when the coding is performed by GPT-4o, with temperature 0.7 and 0, instead of by a human. Thus, the emphasis is no longer on differences across ASR tools but on whether, for a given ASR tool, the coding method produces significant differences. We again use McNemar tests, setting the significance level at 5%. Second, we compute two common measures of IRR: 1) Percentage agreement: proportion of instances in which two “coders” (human versus GPT with default temperature, human versus GPT with temperature 0, and GPT with default temperature versus GPT with temperature 0) produce the same coding outcome. 2) Cohen’s Kappa: a more robust IRR measure that adjusts for agreement

---

<sup>15</sup> We also computed the same number and percentage of different words excluding words that have no strong meaning (option “stopwords = TRUE” of the R library stopwords). The overall patterns remain the same (see SOM6).

by chance, taking values between -1 and 1 (negative values indicate an agreement worse than chance, 0 reflects chance-level agreement, and positive values indicate agreement better than chance).

## 5. Results

### 5.1 Performance to the ASR tools (RQ1)

Table 2 presents the results for the indicators in Block 1. The “Measure” columns present the percentage or average per group, while the “Differences” columns show the differences between pairs of ASR tools for each indicator, including the results of significance tests for these differences.

Table 2. Comparing the three ASR tools (Block 1, human coding)

		Measure			Differences		
		Google	Whisper	Vosk	Google-Whisper	Whisper-Vosk	Vosk-Google
Transcription provided	% Provided	74.0	99.8	90.0	-25.7*	9.8*	15.9*
Length	# Characters	285.0	309.9	385.6	-24.9	-75.7*	100.6*
	# Words	51.1	54.9	68.7	-3.8	-13.8*	17.6*
	# Sentences	1.4	3.3	2.3	-1.8*	0.9*	0.9*
Clarity	% Clear	4.7	21.8	71.9	-17.1*	-50.1*	67.2*
	% Very clear	94.2	72.7	10.7	21.5*	62.0*	-83.4*
Problems	% Missing	8.6	19.2	13.7	-10.5*	5.4*	5.1*
	% Added	6.6	28.5	13.2	-21.9*	15.3*	6.6*
	% Wrong	10.5	5.3	90.3	5.3*	-85.0*	79.8*
	% No problem	76.9	57.4	7.8	19.5*	49.6*	-69.1*
Valid	% Valid	96.7	79.6	93.8	17.1*	-14.2*	-2.9*

*Note.* Total number of observations for “Transcription provided” is 859. The other indicators are only coded for those respondents where an answer was available (N = 636 for Google, N = 857 for Whisper, and N = 773 for Vosk). # stands for average numbers. \* p < 0.05.

These results indicate a potential issue with transcription coverage for Google. Of the 859 respondents for whom we received transcriptions from at least one of the ASR tools, only 74.0% have a transcription from Google. SOM7 compares the 142 voice recordings that lack a transcription in Google but have transcriptions in both Whisper and Vosk to the full set of transcriptions, to learn more about their characteristics. This comparison shows that the files with missing Google transcriptions are longer than average and exhibit more missing words in Whisper and more added words in Vosk, while remaining similar in terms of clarity and valid responses. Thus, the main issue might be related to how Google handles longer voice recordings, and the various parameters that need to be configured in Google.

When Google provides a transcription, however, it performs strongly in terms of clarity, achieving 94.2% of very clear transcriptions and demonstrating low levels of detected missing, added, or incorrect words. Overall, 76.9% of Google transcriptions showed no problems, and 96.7% of its transcriptions were considered valid answers.

Whisper, in contrast, provides transcriptions for 99.8% of these respondents, but has high levels of missing (19.2%) and added words (28.5%), and low rates of valid answers (79.6%). These results support the notion of hallucinations, which inflate answer rates while reducing valid answers. Even when focusing solely on cases with transcriptions available across all three tools – thereby excluding fully hallucinated transcriptions – the number of added words remains significantly higher than for the other ASR tools (19.9%, see SOM5). Additionally, while the percentage of valid answers increases, it still remains significantly lower than for the other tools (88.4%).

Vosk displays an intermediate level of transcriptions (90.0%), likely closer to the percentage of recordings containing intelligible audio (true answer rate). Among respondents with a transcription, Vosk produces the longest outputs in terms of characters and words. However, due to punctuation issues (with Vosk providing almost no periods or commas), this increased length does not translate into a higher sentence count. These punctuation issues contribute to its low clarity, with only 10.7% of transcriptions rated as very clear. This is also associated with a very high rate of incorrect words (90.3% of the transcriptions). For instance, since the question asked about the amount of information provided by nursing homes, many respondents used the word “*información*” (Spanish for “information”) in their answers. However, Vosk frequently mis-transcribed this as “*en formación*” (Spanish for “in training”). Although the mistake was clear, it disrupted the processing of the transcriptions.

When comparing pairs of ASR tools, fewer significant differences are observed between Google and Whisper compared to other pairs. Nonetheless, Google and Whisper still differ notably with respect to most indicators, including clarity and answer validity. Whisper and Vosk as well as Vosk and Google show significant differences across all indicators. The size of these differences is usually large.

Next, Table 3 presents the results for Block 2. Now, the “measure” columns indicate the value of each measure for each pair of ASR tools (e.g., percentage of partly similar meaning between Google and Whisper), while the “differences” columns show differences between pairs (e.g., percentage of partly similar meaning between Google and Whisper minus percentage of partly similar meaning between Whisper and Vosk).

Table 3. Comparing pairs of tools (Block 2, human coding)

		Measure			Differences		
		Google-Whisper	Whisper-Vosk	Vosk-Google	GW-WV	WV-VG	VG-GW
Different words	# Diff. words	17.5	37.4	23.4	-19.9*	14.0*	5.9*
	# Percentage diff. words out of total	18.7	33.0	24.2	-14.3*	8.8*	5.5*
Similar meaning	% Partly similar	10.3	25.8	22.4	-15.5*	3.4	12.1*
	% Very similar	77.1	59.4	73.0	17.7*	-13.6*	-4.1

Note. Overall, we have N = 629 respondents for Google-Whisper (GW) as well as for Vosk-Google (VG) and N = 771 for Whisper-Vosk (WV). # stands for average numbers. \* p < 0.05.

Regarding the word differences across transcriptions, Google and Whisper have the smallest discrepancies, while Whisper and Vosk exhibit the largest ones, both when considering the number of different words and the percentage they represent out of the total words. Similarly, Google and Whisper achieve the highest level of meaning similarity across transcriptions (77.1%), whereas Whisper and Vosk show the lowest one (59.4%).

## **5.2 Comparisons of human and GPT-4o coding (RQ2)**

To compare human and GPT-4o coding, first, we replicate the results from *RQ1*, excluding variables directly created in R, and test for differences between pairs of coding methods: Human versus GPT with temperature 0 (Human-GPT0), Human versus GPT with default temperature (Human-GPT), and GPT with temperature 0 versus GPT with default temperature (GPT0-GPT). Tables 4 and 5 present, respectively, the results for Block 1 and Block 2.

Table 4 shows that only a few significant differences exist between the GPT codings, and that these differences are only observed for Google transcriptions. Thus, setting the temperature to 0 does not substantially affect the results. In contrast, significant and substantial differences are found between human and GPT codings across all three ASR tools. All indicators are significantly different from both GPT codings for Google, six out of seven for Vosk, and four out of seven for Whisper. However, in some cases, these differences do not alter the order of performance among ASR tools. For instance, in terms of validity, Google shows the highest level, followed by Vosk, and then Whisper, both using human or GPT codings (both temperature settings). Nevertheless, in other cases, the order shifts when using GPT instead of human coding. For instance, human coding showed the highest percentage of transcriptions without any of the considered problems for Google, while both GPT codings showed that Whisper has the highest percentage of “no problem.” In general, GPT coding suggests lower quality of the transcriptions. Table 5 also shows significant differences between human and GPT coding, but not between the two GPT codings.

Additionally, we analyze the IRR by computing the percentage agreement and Cohen’s Kappa, for each pair of coding methods. Table 6 presents the results for Block 1 and Table 7 for Block 2.

Again, we observe high similarity between the two sets of GPT codes. The percentage agreement is consistently high across all indicators (minimum = 87.1%). Cohen’s Kappa also indicates strong agreement (minimum = 0.76). In contrast, human and GPT codings exhibit more variability. Depending on the indicator, the percentages of agreement range from 56.3% to 87.6% for Google, 70.9% to 87.0% for Whisper, and 16.0% to 93.8% for Vosk. Cohen’s Kappa is generally low (below 0.20), though there are a few exceptions (maximum = 0.63 when coding validity for Whisper).

## **6. Conclusions and discussion**

### **6.1 Summary**

Our study examined the performance of three ASR tools to transcribe voice answers obtained in a web survey and compared human and GPT-4o coding of the transcribed answers.

Regarding ASR performance (*RQ1*), we found notable differences among the ASR tools. In contrast to Hühne et al. (2025), who found that Whisper more often produces higher quality transcriptions than Google, our results suggest that Google provided the clearest transcriptions

Table 4. Comparing Human and GPT-4o coding with temperature 0 or 0.7, respectively (Block 1)

		Google			Whisper			Vosk		
		Human	GPT0	GPT	Human	GPT0	GPT	Human	GPT0	GPT
Clarity	% Clear	4.7 <sup>ab</sup>	31.4 <sup>a</sup>	30.5 <sup>b</sup>	21.8	23.7	22.2	71.9 <sup>ab</sup>	58.7 <sup>a</sup>	58.1 <sup>b</sup>
	% Very clear	94.2 <sup>ab</sup>	65.9 <sup>a</sup>	66.2 <sup>b</sup>	72.7 <sup>ab</sup>	68.1 <sup>a</sup>	69.3 <sup>b</sup>	10.7 <sup>ab</sup>	7.1 <sup>a</sup>	7.5 <sup>b</sup>
Problems	% Missing	8.6 <sup>ab</sup>	40.7 <sup>a</sup>	42.9 <sup>b</sup>	19.2 <sup>ab</sup>	28.1 <sup>a</sup>	27.3 <sup>b</sup>	13.7 <sup>ab</sup>	97.7 <sup>a</sup>	97.2 <sup>b</sup>
	% Added	6.6 <sup>ab</sup>	36.0 <sup>ac</sup>	38.2 <sup>bc</sup>	28.5	27.9	28.6	13.2 <sup>ab</sup>	78.5 <sup>a</sup>	78.1 <sup>b</sup>
	% Wrong	10.5 <sup>ab</sup>	33.2 <sup>a</sup>	33.6 <sup>b</sup>	5.3 <sup>ab</sup>	11.4 <sup>a</sup>	11.2 <sup>b</sup>	90.3	91.2	91.2
	% No problem	76.9 <sup>ab</sup>	43.9 <sup>ac</sup>	40.7 <sup>bc</sup>	57.4	55.9	55.5	7.8 <sup>ab</sup>	1.7 <sup>a</sup>	2.1 <sup>b</sup>
Valid	% Valid	96.7 <sup>ab</sup>	84.6 <sup>a</sup>	84.3 <sup>b</sup>	79.6 <sup>ab</sup>	65.6 <sup>a</sup>	66.0 <sup>b</sup>	93.8 <sup>ab</sup>	74.6 <sup>a</sup>	74.5 <sup>b</sup>

Note. Superscripts <sup>a</sup>, <sup>b</sup>, and <sup>c</sup> indicate significant differences ( $p < 0.05$ ) between the two columns in which they appear.

Table 5. Comparing Human and GPT-4o coding with temperature 0 or 0.7, respectively (Block 2)

		Google – Whisper			Whisper – Vosk			Vosk – Google		
		Human	GPT0	GPT	Human	GPT0	GPT	Human	GPT0	GPT
Similar meanings	% Partly similar	10.3 <sup>ab</sup>	5.8 <sup>a</sup>	6.8 <sup>b</sup>	25.8 <sup>ab</sup>	9.2 <sup>a</sup>	11.0 <sup>b</sup>	22.4 <sup>ab</sup>	11.9 <sup>a</sup>	12.9 <sup>b</sup>
	% Very similar	77.1	75.0	74.8	59.4 <sup>ab</sup>	65.1 <sup>a</sup>	63.7 <sup>b</sup>	73.0	75.7	75.0

Note. Superscripts <sup>a</sup>, <sup>b</sup>, and <sup>c</sup> indicate significant differences ( $p < 0.05$ ) between the two columns in which they appear.

Table 6. % Agreement and Cohen’s Kappa (Block 1)

	Google			Whisper			Vosk		
	Human-GPT0	Human-GPT	GPT0-GPT	Human-GPT0	Human-GPT	GPT0-GPT	Human-GPT0	Human-GPT	GPT0-GPT
% Agreement									
Clarity	69.2	68.9	90.6	75.4	75.5	95.0	63.5	62.4	87.1
Missing	64.2	62.3	92.8	70.9	72.0	91.0	16.0	16.6	99.2
Added	65.9	64.0	95.0	76.3	76.8	95.6	33.6	33.8	96.0
Wrong	71.1	70.0	95.1	87.0	87.3	97.2	91.6	92.1	99.0
No problem	57.9	56.3	94.3	74.4	73.1	94.8	93.4	93.8	99.4
Valid	87.6	87.3	95.0	84.4	84.8	96.0	80.1	79.7	94.7
Cohen’s Kappa									
Clarity	0	0	0.80	0	0	0.89	0	0	0.76
Missing	0.15	0.15	0.85	0.20	0.22	0.78	0.01	0.01	0.85
Added	0.10	0.09	0.89	0.41	0.43	0.89	0.07	0.06	0.88
Wrong	0.21	0.19	0.89	0.16	0.17	0.86	0.50	0.53	0.94
No problem	0.21	0.21	0.88	0.48	0.45	0.89	0.28	0.35	0.82
Valid	0.30	0.29	0.81	0.62	0.63	0.91	0.30	0.29	0.86

Table 7. % Agreement and Cohen’s Kappa (Block 2)

	Google			Whisper			Vosk		
	Human-GPT0	Human-GPT	GPT0 – GPT	Human-GPT0	Human-GPT	GPT0 – GPT	Human-GPT0	Human-GPT	GPT0 – GPT
% Agreement									
Similar meanings	85.7	86.2	95.3	72.2	72.4	90.4	74.2	73.3	90.3
Cohen’s Kappa									
Similar meanings	0	0	0.88	0	0	0.81	0	0	0.76



with high rates of valid answers, but failed to transcribe a large number of audio files, especially longer ones. Whisper generated transcriptions for almost all cases but exhibited high levels of added words, supporting concerns about hallucinations. Finally, Vosk offered an intermediate transcription rate but suffered from punctuation issues and a high rate of incorrect words, which reduced the clarity of the transcriptions. Pairwise comparisons showed that Google and Whisper have the most similar outputs, while Whisper and Vosk have the largest discrepancies.

Regarding coding methods (*RQ2*), we observed strong agreement between the two GPT coding versions, regardless of the temperature setting. This suggests a limited impact of temperature adjustments on coding. However, significant differences emerged between human and GPT codings across all ASR tools, altering performance rankings in some cases. While human coding ranked Google highest for transcription clarity, absence of problems, and validity, GPT coding ranked Whisper highest for clarity and absence of problems. Additionally, IRR measures were high between both GPT coding versions, whereas human and GPT coding showed substantial disagreement.

## **6.2 Limitations**

This study presents several limitations. Specifically, as we did not have access to the original audio files, we lack a “true value” for comparison. Consequently, we can identify certain transcription issues that impact sentence comprehension, but we cannot fully assess the performance of the ASR tools.

Regarding the coding analysis, we again faced the challenge of not having a true value, as we assess subjective aspects, such as clarity or validity. While we can conclude that GPT and human codings differ, we cannot definitively determine which one is superior.

Finally, both ASR tools and LLMs are evolving rapidly. Thus, the results could change in the near future. The outcomes may also vary when applied to different languages, topics, or target populations. Therefore, further research addressing these aspects is essential to gain a better understanding of the performance of these tools. Other LLMs, such as Gemini Pro or DeepSeek, could also be considered.

## **6.3 Practical implementation**

These findings suggest that each ASR tool has distinct merits and limits. Whisper produces hallucinations (false transcriptions), Vosk has clarity issues and high rates of incorrect words, and Google sometimes fails to provide transcriptions. This seems to be especially the case for longer voice recordings and might be solved by changing the configuration case by case. However, this requires additional time and resources. While Vosk generally performs worse, it aids comprehension in specific instances. Thus, we recommend transcribing the audio files using all three ASR tools. Then, for coding the responses, we propose the following: start with Google’s transcription as the primary option due to its high clarity and validity. If transcriptions are unavailable or unclear, use Whisper as secondary option. Resort to Vosk only when both Google and Whisper transcriptions are absent or unclear. If all three ASR tools produce poor transcriptions, discard the answer.

Regarding the coding method, human and GPT codings differ significantly. Therefore, for now, we recommend continuing with human coding for complex tasks until more evidence on the performance of GPT coding is available. However, given the reduced cost, time, and

effort involved in GPT coding, we also suggest implementing it alongside human coding to identify potential issues in the human coding. Additionally, researchers should carefully design the prompts used for GPT coding. We recommend testing different prompts on small subsets of the data before selecting the one to use for the full dataset. For now, it is also advisable to send the data one answer at a time, as GPT may fail to return a code for some transcriptions if multiple answers are sent at once, though this may improve in the future. Finally, GPT-4o provides slightly different results with each prompt, even when the temperature parameter is set to 0. As a result, it may be beneficial to code each answer at least twice with GPT, and manually review any discrepancies between the two outputs.

However, since the performance of both ASR tools and GPT models is evolving rapidly and can vary with factors like language and background noise for the ASR tools, and specific prompts or settings for GPT, definitive conclusions about which tool is best cannot be made. Researchers should be prepared to adjust to these changes and potentially incorporate newly emerging tools.

Nevertheless, our study provides novel insights that can guide researchers, even as ASR tools and GPT models continue to evolve. These insights remain relevant beyond the specific results presented here and should help navigating future developments. First, the findings highlight the importance of ASR tool selection when transcribing voice responses from web surveys. Researchers should not overlook this decision, as it strongly affects both whether a transcription is generated and its overall accuracy. Second, the recommended approach of using multiple transcriptions to mitigate each ASR tool's weaknesses is likely to remain effective across different contexts and over time. Third, the importance of considering and testing different settings for the tools is also expected to persist. Relying on a single tool with default settings may be a risky strategy. Researchers should test the robustness of their results across tools and settings. Finally, since researchers must make numerous decisions that can affect their results, documenting the choices made is important to ensure transparency and replicability.

## References

- Arlinghaus, C. S., Wulff, C., and Maier, G. W. (2024), "Inductive coding with ChatGPT: An evaluation of different GPT models clustering qualitative data into categories," OSF Preprints, doi:10.31219/osf.io/gpnye.
- Barde, B. V., and Bainwad, A. M. (2017), "An overview of topic modeling methods and tools," 2017 International Conference on Intelligent Computing and Control Systems (ICICCS), doi:10.1109/ICCONS.2017.8250563.
- Callegaro, M., Lozar Manfreda, K., and Vehovar, V. (2015), Web survey methodology. Sage Publishing, <https://study.sagepub.com/web-survey-methodology>.
- Chen, X., Luo, K., Gee, T., and Nejati, M. (2024), "Does ChatGPT and Whisper make humanoid robots more relatable?" In Proceedings in Australasian Conference on Robotics and Automation (ACRA 2023), doi:10.48550/arXiv.2402.07095.
- Conrad, F. G., Keusch, F., and Schober, M. F. (2021), "New data in social and behavioral research," *Public Opinion Quarterly*, 85(S1), 253–263, doi:10.1093/poq/nfab027.
- Deloitte. (2018), "2018 global mobile consumer survey: US edition: A new era in mobile continues." Retrieved from <https://www2.deloitte.com/us/en/pages/technology-media-and-telecommunications/articles/globalmobile-consumer-survey-us-edition.html>.
- Fuller, K. A., Morbitzer, K. A., Zeeman, J. M., Persky, A. M., Savage, A. C., and McLaughlin, J. E. (2024), "Exploring the use of ChatGPT to analyze student course evaluation comments," *BMC Medical Education*, 24, 423, doi:10.1186/s12909-024-05316-2.
- Gavras, K., and 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, doi:10.1080/13645579.2020.1860279.
- Gavras, K., Höhne, J. K., Blom, A., and 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, doi:10.1111/rssa.12807.
- Gummer, T., Höhne, J. K., Rettig, T., Roßmann, J., and 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.
- Gummer, T., Quöß, F., and Roßmann, J. (2019), "Does increasing mobile device coverage reduce heterogeneity in completing web surveys on smartphones?" *Social Science Computer Review*, 37(3), 371–384.
- Gweon, H., and Schonlau, M. (2024), "Automated classification for open-ended questions with BERT," *Journal of Survey Statistics and Methodology*, 12, 493–504, doi:10.1093/jssam/smad015.
- Höhne, J. K., and Claassen, J. (2024), "Examining final comment questions with requests for written and oral answers," *International Journal of Market Research*, 66, 550–558, doi:10.1177/14707853241229329.
- Höhne, J. K., Kern, C., Gavras, K., and Schlosser, S. (2024), "The sound of respondents: Predicting respondents' level of interest in questions with voice data in smartphone surveys," *Quality and Quantity*, 58, 2907–2927, doi:10.1007/s11135-023-01776-8.

- Höhne, J. K., Lenzner, T., and Claassen, J. (2025), “Automatic speech-to-text transcription: Evidence from a smartphone survey with voice answers,” *International Journal of Social Research Methodology*, doi:10.1080/13645579.2024.2443633.
- Höhne, J. K., Revilla, M., and Couper, M. P. (2024), “Pre-registration of the study ‘providing extra incentives for voice answers in web surveys’,” doi:10.17605/OSF.IO/CXZ4S.
- Landesvatter, C., and Bauer, P. C. (2024), “How valid are trust survey measures? New insights from open-ended probing data and supervised machine learning,” *Sociological Methods & Research*, doi:10.1177/00491241241234871.
- Lee, V. V., van der Lubbe, S. C. C., Goh, L. H., and Valderas, J. M. (2024), “Harnessing ChatGPT for thematic analysis: Are we ready?” *Journal of Medical Internet Research*, 26, e54974, doi:10.2196/54974.
- Lenzner, T., Höhne, J. K., and Gavras, K. (2024), “Innovating web probing: Comparing written and oral answers to open-ended probing questions in a smartphone survey,” *Journal of Survey Statistics and Methodology*, 12, 1295–1317, doi:10.1093/jssam/smae031.
- Liu, A., and Sun, M. (2023), “From voices to validity: Leveraging large language models (LLMs) for textual analysis of policy stakeholder interviews,” arXiv, doi:10.48550/arXiv.2312.01202.
- Marion, S. (2024), “How to use OpenAI model temperature?” [Blog post]. Retrieved from <https://gptforwork.com/guides/openai-gpt3-temperature>.
- Mavletova, A. (2013), “Data quality in PC and mobile web surveys,” *Social Science Computer Review*, 31(6), 725–743, doi:10.1177/0894439313485201.
- McMullin, C. (2023), “Transcription and qualitative methods: Implications for third sector research,” *Voluntas*, 34, 140–153, doi:10.1007/s11266-021-00400-3.
- Meitinger, K., van der Sluis, S., and Schonlau, M. (2024), “Keep the noise down: On the performance of automatic speech recognition of voice-recordings in web surveys,” *Survey Practice*, doi:10.29115/SP-2023-0022.
- OpenAI et al. (2024), “GPT-4 technical report,” doi:10.48550/arXiv.2303.08774.
- Pentland, S. J., Fuller, C. M., Spitzley, L. A., and Twitchell, D. P. (2023), “Does accuracy matter? Methodological considerations when using automated speech-to-text for social science research,” *International Journal of Social Research Methodology*, 26, 661–677, doi:10.1080/13645579.2022.2087849.
- Peterson, G., Griffin, J., LaFrance, J., and Li, J. (2017), “Smartphone participation in web surveys,” in *Total survey error in practice*, edited by P. P. Biemer, E. de Leeuw, S. Eckman, B. Edwards, F. Kreuter, L. E. Lyberg, N. C. Tucker, and B. T. West, 203–233. Wiley.
- Proksch, S.-O., Wratil, C., and Wäckerle, J. (2019), “Testing the validity of automatic speech recognition for political text analysis,” *Political Analysis*, 27(3), 339–359, doi:10.1017/pan.2018.62.
- R Core Team. (2023), *R: A language and environment for statistical computing*. R Foundation for Statistical Computing, Vienna, Austria, <https://www.R-project.org/>.
- Radford, A., Kim, J. W., Xu, T., Brockman, G., Mcleavey, C., and Sutskever, I. (2023), “Robust speech recognition via large-scale weak supervision,” Technical report, OpenAI, 2023, <https://cdn.openai.com/papers/whisper.pdf>.
- Revilla, M. (2022), “How to enhance web survey data using metered, geolocation, visual and voice data?” *Survey Research Methods*, 16(1), 1–12, doi:10.18148/srm/2022.v16i1.8013.

- Revilla, M., and Couper, M. P. (2023), “Combining dictation and/or voice recordings with text to answer narrative open-ended survey questions,” ESRA Conference, July 17–21, Milan, Retrieved from <https://www.upf.edu/web/webdataopp/dissemination#Presentations>.
- Revilla, M., and Couper, M. P. (2024), “Exploring respondents’ problems and evaluation in a survey proposing voice inputs,” *methods, data, analyses*, 18(2), 263–280, doi:10.12758/mda.2024.06.
- Revilla, M., and Couper, M. P. (2021), “Improving the use of voice recording in a smartphone survey,” *Social Science Computer Review*, 39(6), 1159–1178, doi:10.1177/0894439319888708.
- Revilla, M., Couper, M. P., Bosch, O. J., and Asensio, M. (2020), “Testing the use of voice input in a smartphone web survey,” *Social Science Computer Review*, 38(2), 207–224, doi:10.1177/0894439318810715.
- Revilla, M., Couper, M. P., and Ochoa, C. (2018), “Giving respondents voice? The feasibility of voice input for mobile web surveys,” *Survey Practice*, 11(2), doi:10.29115/SP-2018-0007.
- Revilla, M., Iglesias, P., Ochoa, C., and Antón, D. (2022), “Webdatavoice: A tool for dictation or recording of voice answers in the frame of web surveys,” OSF, <http://doi.org/10.17605/OSF.IO/B2WYZ>.
- Revilla, M., Toninelli, D., Ochoa, C., and Loewe, G. (2016), “Do online access panels need to adapt surveys for mobile devices?” *Internet Research*, 26(5), 1209–1227, doi:10.1108/IntR-02-2015-0032.
- Roberts, M. E., Stewart, B. M., Tingley, D., Lucas, C., Leder-Luis, J., Gadarian, S. K., Albertson, B., and Rand, D. G. (2014), “Structural topic models for open-ended survey responses,” *American Journal of Political Science*, 58(4), 1064–1082, doi:10.1111/ajps.12103.
- Saldaña, J. (2015), *The coding manual for qualitative researchers*. Sage Publications.
- Schober, M. F., Conrad, F. G., Antoun, C., Ehlen, P., Fail, S., Hupp, A. L., Johnston, M., Vickers, L., Yan, H. Y., and Zhang, C. (2015), “Precision and disclosure in text and voice interviews on smartphones,” *PloS One*, 10(6), Article e0128337, doi:10.1371/journal.pone.0128337.
- Struminskaya, B., Keusch, F., Lugtig, P., and Höhne, J. K. (2020), “Augmenting surveys with data from sensors and apps: Challenges and opportunities,” *Social Science Computer Review*, Advance online publication, doi:10.1177/0894439320979951.
- Theelen, H., Vreuls, J., and Rutten, J. (2024), “Doing research with help from ChatGPT: Promising examples for coding and inter-rater reliability,” *International Journal of Technology in Education*, 7(1), 1–18, doi:10.46328/ijte.537.
- Trabelsi, A., Werey, L., Warichet, S., and Helbert, E. (2024), “Is noise reduction improving open-source ASR transcription engines quality?” In *Proceedings of the 16th International Conference on Agents and Artificial Intelligence – Volume 3: ICAART*, SciTePress, 1221–1228, doi:10.5220/0012457100003636.