

Our Collective Voices: The Social and Technical Values of a Grassroots Chinese Stuttered Speech Dataset

JINGJIN LI*, AImpower.org, USA

QISHENG LI†, AImpower.org, USA

RONG GONG, StammerTalk

LEZHI WANG, StammerTalk

SHAOMEI WU, AImpower.org, USA

The lack of representative stuttered speech data has largely limited the development of stuttering friendly Automatic Speech Recognition (ASR) models. This work studies the first stuttered speech dataset in Mandarin Chinese, created by StammerTalk, a grassroots community of Chinese-speaking people who stutter (PWS). Collected by people who stutter for people who stutter, the dataset includes 50 hours of spontaneous conversations and voice commands from 72 speakers who stutter, capturing stuttered speech with unprecedented scale, diversity, and authenticity. Using this dataset, we are able to benchmark and fine-tune popular ASR models to better understand and mitigate their existing biases against stuttered speech. Our content analysis of the dataset highlights the significant social stigma, overt discrimination, and mental health challenges experienced by PWS in China, exacerbated by the lack of access to scientific knowledge and professional support for stuttering. This dataset not only contributes a critical technical resource for inclusive ASR, but also facilitates self-advocacy and structural changes for PWS in China. By foregrounding lived experiences of PWS in their own voices, we also hope to normalize speech disfluencies and cultivate deeper empathy within the AI community.

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

Additional Key Words and Phrases: AI FATE, datasets, benchmark, speech technology, accessibility, stuttering, stuttered speech

ACM Reference Format:

Jingjin Li, Qisheng Li, Rong Gong, Lezhi Wang, and Shaomei Wu. 2025. Our Collective Voices: The Social and Technical Values of a Grassroots Chinese Stuttered Speech Dataset. In *Proceedings of Proceedings of the ACM Conference on Fairness, Accountability, and Transparency (FAccT '25)*. ACM, New York, NY, USA, 21 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 INTRODUCTION

Trained and optimized for fluent speech, existing automatic speech recognition (ASR) systems performs poorly for people who stutter (PWS): it often cuts them off from speaking and interprets the speech of PWS with a multitude higher error rates than average [3, 21]. As ASR systems become a ubiquitous part of today’s communication ecosystem, their inability to work with stuttered speech not only creates additional barriers for PWS to interact with popular

*Both authors contributed equally to this research.

†Both authors contributed equally to this research.

Authors’ addresses: [Jingjin Li](mailto:jingjin@aimpower.org), jingjin@aimpower.org, AImpower.org, Santa Clara, California, USA; [Qisheng Li](mailto:qishengli@aimpower.org), qishengli@aimpower.org, AImpower.org, Seattle, Washington, USA; [Rong Gong](mailto:rong.gong@stammertalk.net), rong.gong@stammertalk.net, StammerTalk; [Lezhi Wang](mailto:lzwangcn@gmail.com), lzwangcn@gmail.com, StammerTalk; [Shaomei Wu](mailto:shaomei@aimpower.org), shaomei@aimpower.org, AImpower.org, Mountain View, California, USA.

Permission to make digital or hard copies of part or all 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 third-party components of this work must be honored. For all other uses, contact the owner/author(s).

© 2025 Copyright held by the owner/author(s).

Manuscript submitted to ACM

products and services like smart speakers, in-car navigation systems, and automatic phone menus, but also leads to psychological harms and socioeconomic disadvantages [8, 23, 43, 45].

The lack of adequate stuttered speech data has been a key bottleneck in addressing ASR performance disparities [22, 26]. Existing stuttered speech datasets, such as FluencyBank [32] and LibriStutter [19], primarily focus on English speech and are often limited in size, representativeness, and annotation consistency [22]. Recent efforts in collecting atypical speech for inclusive ASR also face challenges in engaging and authentically representing the stuttering community [26]. Distinct from other speech etiologies such as Amyotrophic Lateral Sclerosis (ALS), stuttering is highly variable and inherently social: the severity and patterns of stuttering can vary significantly across individuals, situations, and conversation partners [9, 38]. Another major shortcoming of these datasets is their lack of unscripted conversational data, which is crucial for use cases like auto captioning and conversational agents. Typically, conventional speech data collection approaches, which record participants' monologue responses to given prompts [2, 26, 28], work poorly at capturing authentic, real-world stuttering patterns.

To fill these gaps, StammerTalk - a grassroots community of Chinese speaking people who stutter - self mobilized to create the first and largest stuttered speech dataset in Mandarin Chinese [24]. Containing about 50 hours of **Mandarin speech with both spontaneous conversations and voice command dictations** from 72 individuals who stutter, the StammerTalk dataset offers unprecedented versatility and authenticity compared to existing stuttered speech datasets. Through statistical and content analysis of the dataset, as well as the evaluation and fine tuning of state-of-the-art ASR models using this data, our study shows the unique value of community-created stuttered speech datasets in capturing the heterogeneity and variability of stuttering, highlighting its efficacy in uncovering fairness issues in existing ASR models and raising awareness of the needs of PWS among designers and developers of ASR models. The content analysis of the dataset highlights the significant social stigma, overt discrimination, mental health challenges, experienced by PWS in China, and the lack of access to scientific knowledge of and professional support for stuttering. This dataset holds significant social and educational value, serving as a platform for self-advocacy and public discourse on stuttering in China, and cultivating deeper empathy within the AI community. Our work contributes to existing literature on ASR fairness, inclusivity, and equity, advocating for open and respectful collaboration between the research community and disability communities in the full cycle of speech AI design and development. Curating the dataset also exemplifies justice and empowerment by placing the agency of PWS at the heart of its creation. By being curated by PWS, with PWS, and for PWS, the process empowers the community to take ownership of the narratives and resources that shape their representation in AI, challenging the traditional top-down approaches often seen in AI research.

2 RELATED WORK

2.1 Stuttering

Stuttering is a neurodevelopmental condition that impacts PWS in behavioral, emotional, and cognitive aspects [4]. Besides speech struggles, PWS often develop adverse emotional and cognitive reactions of stuttering, such as fear, guilt, shame, and self censorship. The emotional, and cognitive components of stuttering make it highly variable and diverse [4, 9, 38]. Across individuals who stutter, a wide range of speech characteristics exist, from observable speech disfluencies like word repetitions, to "masked" disfluencies such as word switching and circumlocution [9]. Even for the same speaker, their stuttering pattern and frequency can vary from complete fluency when speaking to oneself alone, to severe disfluencies when speaking over the phone or to a group [38]. Such inherent irregularities in stuttered speech make it particularly challenging when applying modeling techniques that rely on statistical methods

and pattern recognitions – a dominant approach in modern ASR. As a result, today’s ASR systems have great difficulty in understanding and interacting with users who stutter [3], showing as high as 50% WER for severely stuttered speech, 10 times the reported consumer average of 5% [21]. Our work contributes to a deeper understanding of AI FATE issues for people who stutter, a population profoundly affected yet often overlooked by speech AI technologies.

Despite prevalent stigma [7] and discrimination [6, 8, 11, 15], the stuttering community has been actively pushing back on ableist expectations of speech fluency and advocating for broader acceptance and respect of stuttered speech in everyday communications [9, 10, 36]. While previous survey study showed PWS in China have more adverse experiences related to stuttering compared to PWS in western or developed countries [25], the creation of StammerTalk dataset not only contributes intimate, personal accounts of the lived experiences of stuttering in China, but also facilitates self advocacy and structural changes for the Chinese stuttering community.

2.2 Inclusive ASR for Stuttering

Ableist assumptions about the syntax and temporality of human speech are prevalent in the design and development of ASR. For example, endpointer models – a component of ASR that identifies the end of an utterance – frequently truncate the speech input from PWS as these models have learned auditory features, including the duration of silence between sounds, from fluent speech [21]. Similarly, ASR decoders tend to inject unrelated words in place of partial or whole word repetitions [21, 29]. Recent technical explorations show promise in addressing these biases by fine-tuning ASR decoders with disfluent speech [21, 29], adjusting endpointer thresholds [21], and model personalization [14, 40]. In parallel, research on stuttering event detection [1, 19] explores a two-step “detect and pass” approach to help ASR models better process stuttered speech [34].

However, current solutions often require a relatively large amount of annotated personal speech data [14, 40], with improvements limited to short phrases [14], prompted speech [14], pre-defined vocabularies and scenarios [21, 29, 40], and speech with reduced articulations (e.g. speech by patients with ALS, CP, or Parkinson’s Disease) [14, 40]. The question remains regarding their generalizability to everyday interactions with ASR systems by PWS “in the wild”.

Even when ASR systems do manage to transcribe stuttered speech, the resultant transcriptions often remove disfluencies like filler words, inevitably reinforcing the ableist construct of stuttering as “errors” and “undesirable” [36] and denying the opportunity for PWS to have their speech transcribed authentically [24].

Overall, there is still a significant gap in speech technology accessibility for PWS. Our work aims to address it by auditing popular ASR models using stuttered speech across diverse scenarios and speaker profiles.

2.3 Stuttered Speech Datasets

To address the lack of diverse speech data for inclusive and robust ASR models [27, 29, 48], there have been several industry initiatives to collect diverse speech samples across languages [2], accents [2, 28], and speech disabilities [26, 42]. However, despite these efforts, high-quality stuttered speech data remain scarce in the public domain.

Considered highly sensitive personal data, stuttered speech collected and annotated by companies is often inaccessible to the broader research community [21, 26, 29, 34, 48]. Datasets collected by academic researchers, such as FluencyBank [32] and UCLASS[16], were limited in sizes and annotation consistency, primarily used for speech therapy [22]. Open and scalable datasets, such as SEP-28k [22] and LibriStutter[19], also have various shortcomings in terms of annotation completeness, representativeness, and authenticity. Collected from public podcasts by people who stutter, SEP-28k consists of 28K 3-second audio clips labeled with only stuttering events (i.e. whether the clip contains prolongation, sound repetition, etc.) but lacks text transcriptions [22]. While LibriStutter does provide transcriptions, it

contains no real speech from PWS but synthetic stuttering utterances (e.g. repetitions, prolongations, interjections) injected into audio books read by fluent speakers [19]. By introducing StammerTalk dataset, the first and largest corpus of stuttered speech in Mandarin Chinese, our work advances efforts to represent stuttered speech and its community in AI, and underscores the importance of community-driven data collection to authentically capture stuttering’s diversity and amplify the community’s voice in AI development.

3 METHOD

The dataset was created by StammerTalk (口吃说) community¹, an online, grassroots community of Chinese speaking people who stutter. Speech data collection was conducted by two StammerTalk volunteers, who also stutter, with participants who stutter over videoconferencing platforms. The recorded speech contains both unscripted conversations between the volunteer and the participant, and the dictation of a list of 200 voice commands by the participant. 70 adults who stutter (AWS) participated in the recording with two StammerTalk volunteers, resulting in a dataset of 50 hours speech from 72 AWS. The recorded speech was transcribed semantically and verbatim, with five distinct stuttering event annotations ([] - Word-level repetition, /r - sound repetition, /b - blocks, /p - prolongation, /i - interjection) embedded in markups. Obtaining verbatim transcription that includes word repetitions (e.g. “My, my, my name”) and interjections (e.g. “hmm”) was a deliberate choice made by the StammerTalk community, to allow disfluencies respected and preserved by ASR models rather than being automatically removed. The annotation was performed by professional speech data annotators, and reviewed by a StammerTalk volunteer. More details about the data collection and annotation process can be found in previous work [24].

To understand the characteristics and quality of the StammerTalk dataset, we perform both quantitative and qualitative analysis on its technical and social properties.

3.1 Quantitative analysis

We first conduct descriptive analysis of the StammerTalk dataset, comparing its scale and speech diversity with existing stuttered speech datasets. We also benchmark the performance of prominent ASR models with our dataset to measure and diagnose ASR biases towards stuttered speech. Lastly, we explore ASR model fine tuning using the StammerTalk dataset and demonstrate the value of data diversity in improving the fairness of foundational models.

3.1.1 Descriptive analysis. Stuttering is not a monolith. The frequency and types of stuttering can vary significantly across individuals and situations - a common source of insecurity and frustration for PWS [38]. While existing stuttered speech datasets often suffer from limited scale and coverage of the heterogeneity within stuttering [21], we measure the scale and speech diversity of the StammerTalk dataset in terms of speakers, speaking tasks, stuttering frequency and severity, and speech variability among and within speakers.

3.1.2 Benchmarking. To understand ASR’s ability to transcribe and respect speech disfluencies, we audit two state-of-the-art ASR services – Whisper (large-v3)² and wav2vec2.0 (large-chinese-zh-cn)³ – with two types of ground truth transcriptions: 1) a **semantic** transcription with word repetitions and interjections excluded; 2) a **literal** transcription with the stuttered utterances kept verbatim. We remove all stuttering event markups in both cases.

¹<https://www.stammertalk.net/>

²<https://github.com/openai/whisper>

³A fine-tuned version of wav2vec2.0 optimized for Mandarin speech, see <https://huggingface.co/wbbbbb/wav2vec2-large-chinese-zh-cn>

We calculate the character error rate (CER), a metric commonly used to measure the ASR performance in Mandarin Chinese, using both semantic and literal transcriptions as references. CER measures the errors in model generated transcriptions at the character level, including substitutions (SUB), insertions (INS), and deletions (DEL).

3.1.3 Model fine tuning. To further evaluate the technical value of the StammerTalk dataset, we fine-tuned the LoRA adapter for the Whisper-v2-large model [46] on the StammerTalk Conversation dataset using literal transcriptions as references. The dataset was divided into a train/dev/test split, ensuring a balanced representation of mild, moderate, and severe stuttering levels in each subset. Specifically, 65% of the dataset was allocated for training, 10% for dev, and 25% for testing. This split strategy ensured robust evaluation of the model's performance across all severity levels.

Fine-tuning was performed with a training objective to minimize the character-level transcription error, focusing on preserving stuttered speech patterns such as word repetitions and interjections. The model was fine-tuned using 3 epochs, with early stopping applied based on the validation loss to avoid overfitting. Training hyperparameters included a learning rate of $1e-3$, batch size of 16, and the AdamW optimizer.

The fine-tuned model's performance was evaluated on the held-out test set using the same character error rate (CER) metrics as in the benchmarking task. The CER was further analyzed by severity level (mild, moderate, severe) to assess how well the fine-tuned model handles varying degrees of disfluencies compared to the baseline Whisper model.

3.2 Qualitative analysis

The StammerTalk dataset is unique as it contains of 70 spontaneous conversations between two people who stutter [24]. While the conversations were unscripted, most of them naturally converged on shared experiences and personal stories around stuttering [24], making the StammerTalk dataset the first public archive of lived experiences of PWS in China to the best of our knowledge. To unpack the collective narratives captured in the StammerTalk dataset, we used an inductive open-coding analysis approach [33] to conduct the content analysis of recorded conversations. Our qualitative analysis consists of the following steps:

- (1) First, the first two authors and the last author independently reviewed the transcripts of the first five participants and generated initial codes by adding comments directly to the documents. For example, we had comments "Feeling ashamed after stuttering during meetings" to describe emotional feelings after stuttering.
- (2) The three researchers then met to read through and discuss the transcripts together. Through this discussion, they refined their initial comments into a set of agreed-upon codes, organized these codes into broader categories, and developed a preliminary coding scheme. For example, codes like "Stuttering is from imitation", "Stuttering can be cured" were grouped under the category "Misconception of stuttering".
- (3) The first author then thoroughly reviewed the remaining transcripts multiple times, applying codes as comments and continuously refining the coding scheme through an iterative process. See coding scheme in appendix C.
- (4) In subsequent research meetings, the team collaboratively identified key thematic insights emerging from the categorized codes and synthesized these insights for reporting.

4 FINDINGS

4.1 Descriptive Statistics

Our descriptive analysis of the StammerTalk dataset highlights its scale and data diversity, illustrating its unique quality to represent stuttered speech for ASR.

Table 1. **Dataset scale and scope as characterized by speech duration (Duration), the number and types of speakers (Speakers), whether it provides speech transcription (Transcription), types of speaking tasks (Tasks), and Language.**

Dataset	Duration	Speakers	Transcription	Tasks	Language
FluencyBank* [32]	3.5 hrs	32 AWS	Yes	conversation, reading article	English
LibriStutter [20]	20 hrs	50 non-PWS	Yes**	audiobook	English
UCLASS* [16]	53 mins	25 CWS	Yes	conversation	English
SEP-28k [22]	23 hrs***	not reported	No	podcast	English
StammerTalk	50 hrs	72 AWS	Yes (verbatim)	conversation, voice commands	Chinese

* Limited to the transcribed portion of the dataset.

** Stuttered utterances are masked in the transcription as “STUTTER”.

*** Split into 28K 3-second clips.

Abbreviations: AWS - adults who stutter; CWS - children who stutter.

4.1.1 Scale and Scope. We measure the scale and scope of the StammerTalk dataset in terms of speakers, speech duration, stuttering events, and speech content. Key statistics for these aspects are provided in Table 1, along with existing datasets for comparison.

Speakers and Duration. A total duration of 50-hours speech data were included in StammerTalk dataset from 72 speakers. Excluding the two StammerTalk volunteers, most participants (64 out of 70) are from mainland, China. 34% (24) of the participants are female, much higher than the reported 20% or less among adults who stutter. Each participant contributed on average 33.0 minutes of conversational speech ($min=17.2$, $max=49.9$, $SD=7.32$), with an average of 17.8 minutes ($min=7.14$, $max=34.93$, $SD=5.6$) and an average of 15.23 minutes ($min=6.45$, $max=27.6$, $SD=5.23$) of voice command dictation. Many participants found speaking with another PWS both rare and pleasant [12, 24], thus spent more time on the conversations.

Stuttering Events. A total of 28,310 stuttering events were annotated in the StammerTalk dataset. Table 2 compares the frequency and distribution of annotated stuttering events in the StammerTalk dataset with existing stuttered speech datasets with stuttering event annotation, highlighting the quantity and diversity of stuttering events captured in the StammerTalk data. We also compute the *Average Stuttering Rate* by dividing the total count of stuttering events by the duration of the speech, and find that conversational speech in the StammerTalk dataset exhibits approximately 25% more stuttered utterances compared to the stuttering podcast (SEP-28k) and synthetic stuttered speech (LibriStutter).

Table 2 also shows *Event Type Distribution*, the percentage of each stuttering type among all annotated stuttering events. We note that a direct comparison between the Sep-28k and StammerTalk datasets may not provide the full picture, as StammerTalk’s event annotation is performed at the character level, which offers greater granularity than

Table 2. **Overall frequency and distribution of annotated stuttering events.**

	Avg. Stuttering Rate (per minute)	Total Stuttering Events	Event Type Distribution				
			[]	/b	/p	/r	/i
LibriStutter [19]	12.5*	15,000*	20%	20%	20%	20%	20%
SEP-28k [22]	12.26	17,267	16%	19%	16%	14%	35%
StammerTalk: Conversation	15.83	19,674	42%	6%	18%	9%	25%
StammerTalk: Dictation	8.10	8,636	53%	8%	17%	16%	6%

* Stuttered utterances were synthetically generated.

the clip-level annotation in Sep-28k. However, we do observe a significant shift towards more word and phrases repetitions and less sound repetitions in the StammerTalk dataset, signaling potential phonological differences between stuttering in Chinese and in English. Meanwhile, we notice that SEP-28k dataset contains 40% more interjections than in StammerTalk Conversations, which could be attributed to different definitions of interjections in these two datasets: while SEP-28k considers any filler words - such as “um,” “uh,” and “you know” - as stuttering interjections, StammerTalk’s annotation excludes natural interjections that blend into the speech flow.

Stuttering Transcription. The StammerTalk dataset contains both voice command dictation and unscripted conversations in Chinese. Excluding stuttering event annotations, the verbatim text transcription of StammerTalk dataset contains 425K Chinese characters (274K for conversations, 171K for voice command dictation).

To summarize, the StammerTalk dataset surpasses existing datasets in its duration, speakers, and stuttering frequency. It contains 20 times more transcribed speech data from people who stutter than what is available today (i.e. FluencyBank), and a multiplied number of speakers who stutter. Additionally, it provides both stuttering event annotations and verbatim transcriptions, enabling versatile applications across a wide range of technical domains. Unlike podcasts or audio books, the StammerTalk dataset contains unscripted conversations and voice command dictations that closely resembles real-world speech product use cases, such as meeting transcriptions and speech-operated devices.

4.1.2 Speech Diversity. Contrasting to previous stuttered speech datasets [19, 22], the StammerTalk dataset captures a wide spectrum of stuttering frequency and patterns across PWS in different scenarios, providing a much more comprehensive representation of the variability and heterogeneity of stuttered speech for speech AI.

Stuttering frequency. While all participants self-identified as people who stutter, their frequency of stuttering varies greatly. To quantify individual stuttering frequency, we calculate *disfluency rate*, as defined in [13, 21], by dividing the total number of stuttering events over the total transcribed non-stuttering character count for each speaker. Following conventional approaches [13, 21], we categorize speakers into three groups based on their *disfluency rates*, corresponding to mild (0-5%), moderate (6-20%), and severe (over 20%) stuttering. As illustrated in Fig. 4 (in Appendix A), while the participants stutter more in Conversations (mean=9.2%) than in Command Dictation (mean=7.1%), the stuttering frequency varies more in Command Dictation (std=0.15) than in Conversation (std=0.08). As a result, the grouping of speakers varies across two tasks: for Conversation, 20, 44, and 6 speakers are categorized as mild, moderate, and severe stuttering, respectively, whereas for Dictation, the numbers are 46, 18, and 6.

The variation in *disfluency rates* across different tasks and speakers highlights the dynamic and situated nature of stuttering: its severity varies not just across individuals but also within the same individual. For some, reading is much easier than conversations; whereas for others, reading could be extremely challenging (100% *disfluency rate*).

Stuttering patterns. PWS often stutter differently: some speak with more repetitions, some frequently block, while some stutter covertly [9, 38]. Fig 1 shows the breakdowns of annotated stuttering events, for all 70 participants, highlighting the variation with their stuttering patterns. It also illustrates the change in stuttering patterns for the same speaker with different tasks: participants often have relatively more interjections in conversations, but show increased sound repetitions when dictating commands.

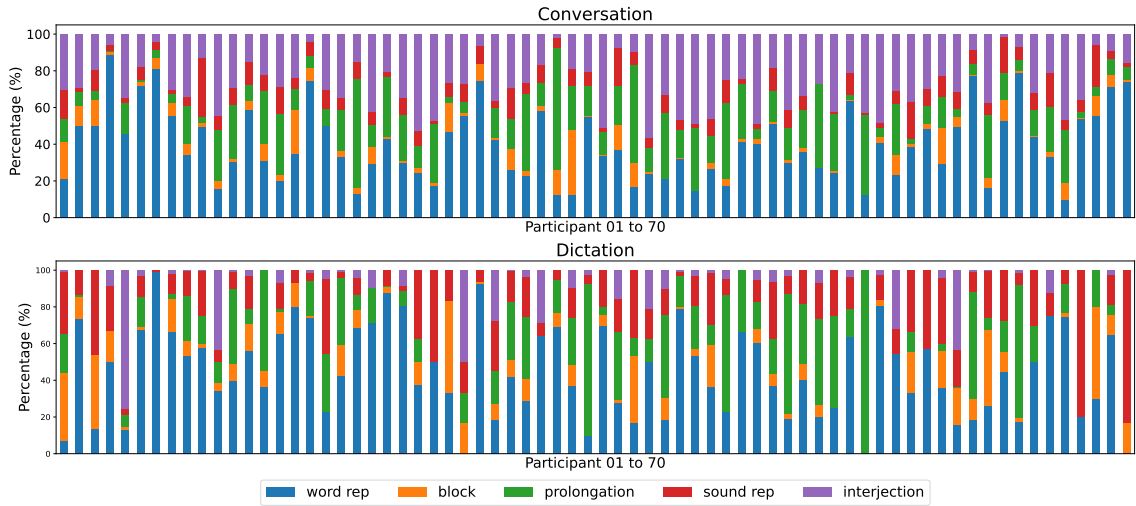


Fig. 1. Breakdowns of five annotated stuttering events for 70 participants

4.2 Benchmarking Results

The benchmarking results with Whisper model are shown in Fig. 2. We notice that it performs reasonably well with semantic transcriptions of mildly stuttered speech, achieving a CER of 10.74% for unscripted conversation (Fig. 2a). However, CER increases with stuttering severity, reaching 13.46% for the moderately and 19.46% for the severely stuttered speech. These CERs are both higher than the reported performance for the general population, namely, 12.8% on the Common Voice 15 dataset and 7.7% on the FLEURS dataset⁴.

Comparing Fig. 2a with Fig. 2b, we find a sharp increase in deletion errors (DEL) when referencing on literal transcriptions. Further inspection of the results shows that Whisper has difficulties in generating disfluent literal transcriptions, often “smoothing” its transcriptions by removing repeated words or phrases. We provide examples of this behavior in Table 3 in Appendix B. As presented in Fig. 2c and Fig. 2d, we find that CERs are higher for Dictation tasks compared to Conversation, potentially due to its reliance on language model to “guess” correct transcription using the semantic context, which is more limited for voice commands. The wav2vec model, in contrast, performs 1.5 to 2 times worse than Whisper, and produces a lot more substitution mistakes. Manual inspection finds that wav2vec model often substitutes a character with its homophones, undervaluing the semantic context. More detailed results for wav2vec model can be found in Appendix B.

4.3 Fine Tuning ASR Models with StammerTalk Dataset

The results of fine-tuning Whisper with the literal transcriptions of StammerTalk Conversation are presented in Fig 3. We observe substantial improvements in transcription accuracy across all severity levels of stuttering compared to the baseline model. For mildly stuttered speech, the fine-tuned model achieves a CER reduction from 16.34% (baseline model) to 5.8%, closing Whisper’s performance gap between stuttered and fluent speech [31]. Similarly, for moderately and severely stuttered speech, the CER drops from 21.72% to 9.03% and from 49.24% to 20.46%, respectively. Our results illustrate the effectiveness of fine-tuning in improving ASR accuracy for stuttered speech across all severity levels.

⁴<https://github.com/openai/whisper?tab=readme-ov-file>

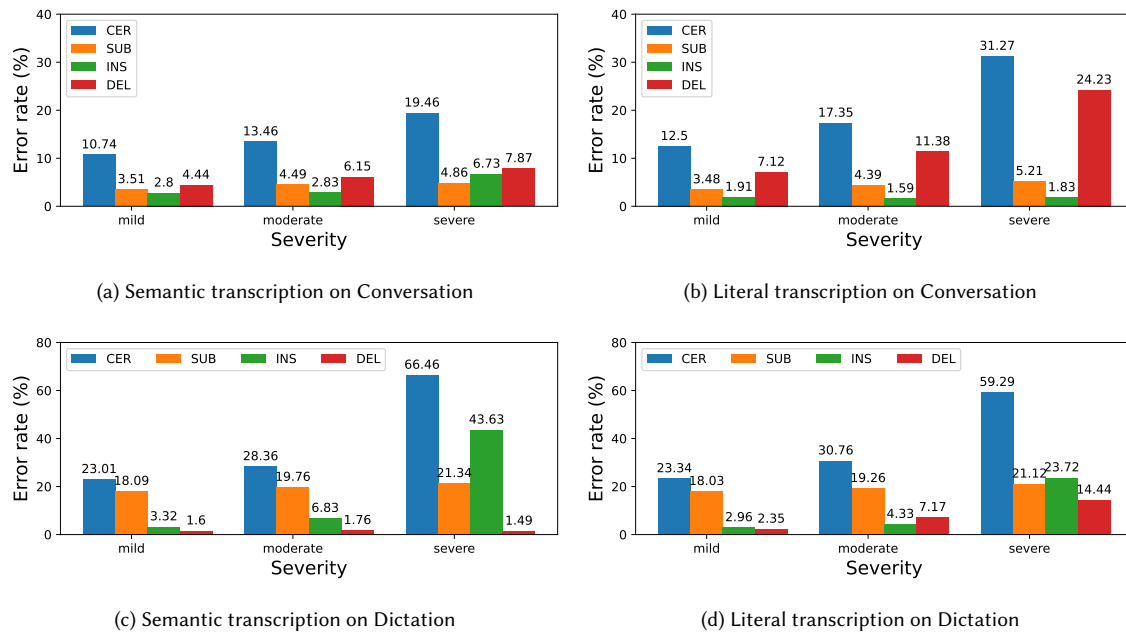


Fig. 2. Character error rate (CER), substitution (SUB), insertion (INS) and deletion (DEL) error rates for Whisper evaluated on StammerTalk Conversation and Dictation.

We want to call out the significant reduction in deletion errors (DEL) after fine-tuning. As shown in Fig 3, the DEL rate drops from 26.56% to 2.29% for severely stuttered speech, and from 15.77% to 1.27% for moderately stuttered speech. Consistent with our benchmarking results, the baseline Whisper-large-v2 model often smoothens its output by omitting repeated words or phrases. This behavior, while generating more fluent transcript, leads to higher deletion error (DEL) rates when evaluated against verbatim transcriptions. We find the fine-tuned model more inclusive of speech disfluencies: it is more likely to preserve disfluencies rather than erasing them from the generated transcript.

Overall, fine-tuning Whisper with the StammerTalk dataset helps the model better recognize and preserve speech disfluencies and significantly improves its transcription accuracy with stuttered speech. Our results demonstrate both the importance and the effectiveness of model fine-tuning with StammerTalk dataset in addressing ASR’s fluency biases.

4.4 Qualitative Findings: Lived Experience of Stuttering in China

Our qualitative analysis find that, the conversations between the participant and the interviewer, although unscripted, have often centered on stuttering and the lived experiences of PWS. This aligns with stuttering’s socially isolating nature [5, 41] and having a one-on-one conversation with another PWS was reported as among a key motivator for joining the study [24]. We report qualitative findings that focus on the lived experience of stuttering in China. Participants’ quotes were translated into English and lightly edited for readability.

4.4.1 Prevalent social stigma and psychological impact. While stuttering is known to lead to negative emotional and cognitive reactions in PWS, the prevalence of social stigma towards stuttering, as its associated strong psychological impact, stands out in our data. All participants report experiencing some form of systematic discrimination and stigma

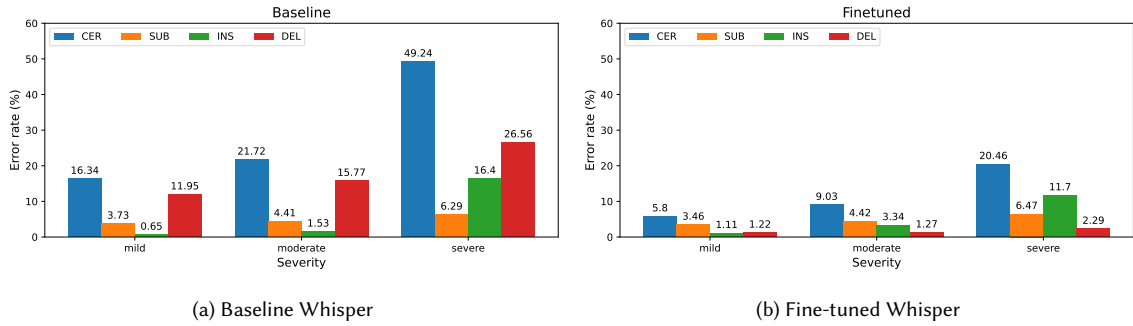


Fig. 3. Character error rate (CER), substitution (SUB), insertion (INS) and deletion (DEL) error rates for baseline and fine-tuned Whisper-large-v2 on StammerTalk Conversation.

towards stuttering at home, school, or workplace environments, which significantly impacted their social interactions and overall well-being. Nearly a third of the participants bring up the experience of being bullied by peers during childhood due to their stutter, which led to lasting psychological trauma, social anxiety, feelings of inferiority, avoidance behaviors, and depression. For example, P20 shared,

"I was afraid of speaking in front of many people. I was also scared of being called on by the teacher to answer questions in class. Sometimes, when I got very nervous, I couldn't speak at all. My classmates would laugh at me, making my childhood feel quite oppressive."

The psychological toll of stuttering was so severe in P27's case that she developed self-harm to punish herself for speech difficulties during childhood,

"When I was a child, every time I got stuck or stumbled while speaking, or when I repeated words, I would severely punish myself. One way I punished myself as a child was by keeping my fingernails very long. If I couldn't speak properly, I would clench my fists, and my fingernails would dig into my skin."

Several participants report a lack of understanding and acceptance of their stuttering by family members, which further undermines their psychological well-being. P69 reports her experience of depression and anxiety because of stuttering combined with misunderstandings from her parents.

"My parents are really very stubborn and not open-minded...I developed depression and anxiety because of my stuttering, but I had no way to tell them about it...and then I took a leave of absence from school. Gradually, my stuttering caused serious emotional problems, which also affected my physical health. I became extremely anxious and developed some psychosomatic symptoms."

For a few participants who report neutral or positive experiences with stuttering, they often attribute their experiences to the acceptance attitude. As P26 reflects the acceptance mindset has minimized the stuttering impact of stuttering on her life, *"What helped me the most was a shift in my mindset. I began to realize that although stuttering does have some impact on my daily life, as long as I handle it properly, its impact on me is actually very small. So now my attitude toward it is to accept it in a healthy way."*

4.4.2 Workplace and professional discrimination. Although the stigma is prevalent, many participants highlight the tension between stuttering and the competitive employment environment in China today. As noted by P12,

"If I weren't working, I would feel that my stuttering isn't particularly severe and that I can adjust it in time. However, in today's society, there are very high demands on a person's overall abilities and competitiveness. If you want to be competent in certain positions, it's important to avoid having any weaknesses. Having a stutter does impact me personally and can also make it difficult to perform well in certain roles."

Career choices are greatly affected by stuttering, pushing many participants into careers with minimal verbal interactions. P30 shares *"I feel quite anxious (about my stuttering) so ... I don't dare to choose a job that requires a lot of speaking. That's why I'm currently doing research work."* While occupational risks and labor discrimination for PWS were also reported in the US [15], they were much more overt and socially accepted in China, according to the StammerTalk dataset participants. In particular, having a stutter could disqualify someone for professional fields such as teaching and healthcare, as PWS are assumed to not be able to meet the verbal communication demands. P27 noted she was discouraged from becoming a medical doctor because *"In the handbook for college applications, it stated that people who stutter are prohibited from applying for clinical majors"*.

While stuttering is highly dynamic and variable across individuals and situations, environmental stressors - such as time pressure and listeners' reaction - are reported to lead to more severe stuttering [38]. Overall, participants report fewer stuttering episodes when communicating with familiar individuals or in intimate settings but greater struggles during presentations, interviews, or interactions with strangers - situations common in the workplace. Stuttering is thus often viewed by employers, and internalized by our participants, as a failure and a sign of incompetence at work. P25 explained, *"in a competitive environment, I don't want to fall far behind my peers. If my supervisor knows that I have a stutter, they might not offer me important opportunities."*

Prevalent workplace discrimination of stuttering, combined with the lack of structural protection from labor unions or employment laws, drives our participants to spend significant efforts to "fix" their stutter or at least "pass" as fluent. For example, despite the discouragement from the college application handbook, P27 applied for medical school - while hiding her stutter - to "help others with similar hardship" and eventually became a physician. Although she was well appreciated by her patients for her skills and patience, she still felt constant pressure to speak fluently and would sometimes take sleeping pills or drink alcohol to reduce stuttering.

4.4.3 Coping mechanism. In response to the stigma and discrimination towards stuttering, participants develop various coping mechanisms, predominantly focused on concealing their stutter, fluency shaping techniques, and avoidance. Avoidance strategies include avoiding certain words or sounds and substituting challenging vocabulary, avoiding speaking situations and relationships. Some participants report avoiding communication as a strategy to manage the fear or reality of stuttering, which may lead to increased feelings of loneliness and isolation. For example, P51 intentionally refrains from joining conversations with colleagues,

"In the office, I rarely initiate conversations with others. Sometimes, even when I'm interested in what they're talking about, I avoid joining in just to prevent stuttering. To some extent, it feels a bit suppressive, but I can accept it because staying silent feels better than stuttering. In a way, it's like closing myself off."

Although around one-third of the participants have sought "professional" help to manage stuttering, including attending online and in-person stuttering correction programs, seeing a pediatrician, a stomatologist (mostly during childhood), or other medical professionals, only a few have visited professionally trained speech language pathologists (SLPs). Some participants share that they have attended various speech programs claiming to cure stuttering but gained little or no improvement in speech fluency afterwards. For example, P28 expressed her disappointment,

"I have attended a stuttering correction class but the experience was very disappointing. Not only did it fail to make my speech more fluent, but it also increased my frustration with myself. The class promoted the idea that if you don't speak fluently, it's entirely your fault — you're not using the methods correctly or not practicing breathing properly."

These programs are often expensive yet not effective, so some participants seek resources online, reading books about stuttering or joining stuttering support groups instead. Participants express that these resources are most helpful in improving their acceptance towards stuttering: *"I started participating in in-person stuttering support groups, and by hearing other people who stutter share their experiences, I was slowly able to accept my stuttering. Even if others outside couldn't accept it, I felt that my mindset had changed."* (P16)

In contrast to the documented benefits of self-disclosure [47], more than half of the participants report they often avoid disclosing their stuttering. Participants fear that disclosing their stuttering could lead to misunderstandings or negative interpersonal and professional consequences. P4 describes the discomfort and the interpersonal risks of disclosing stuttering, *"I'm afraid that if I disclose my stuttering to my friends, they might leave me or dislike me."*

4.4.4 Misconceptions about stuttering. We observed a general lack of scientific understanding about stuttering in our data, even within the stuttering community in China. Such deficit of knowledge perpetuates harmful stereotypes and increases social and self stigma. One common misconception is that many people think they acquired their stuttering from imitating stuttered speech during childhood. Similarly, some PWS also worry that their children might develop a stutter by imitating them, which causes them significant psychological stress, P47 expressed,

"I might unintentionally influence my child, because young children naturally imitate their parents. I feel that my stuttering not only affects me but could also impact my children's future, including their job interviews, career opportunities, and even their romantic relationships."

Participants also report a prevalent view that verbal fluency reflects cognitive competence, which is often used to justify social exclusion, discrimination, and limited opportunities for people who stutter. P8 shared, *"Most people have a misunderstanding about stuttering: they assume that people who stutter also have low intelligence."*

Some participants report that their stuttering has been treated as a physical abnormality rather than a complex neurological and psychological condition. As P34 mentioned *"My parents thought that my speech issue was due to a physiological condition. So they took me to have surgery to shorten my tongue frenulum, but it didn't improve my stuttering."*

The lack of scientific understanding of the causes and nature of stuttering could lead to unrealistic expectations for PWS to speak fluently and harsh criticism towards them for not trying hard enough. As P56 shared,

"My parents would criticize me harshly about my speech if I didn't speak well. They would say that I must speak properly and that if I couldn't, it would be difficult for me to find a good job in the future...For as long as I can remember, whenever I didn't speak well, they would always criticize me."

Despite widespread misconceptions, StammerTalk data collectors – both of whom resided outside China and received more comprehensive stuttering therapy and professional support – frequently shared information with participants on various aspects of stuttering, including its cause, techniques to improve fluency, and available resources for management. Thus, the data collection process also served as an educational opportunity for PWS to learn about stuttering and reflect on their personal experiences. For instance, P31 commended the interviewer for sharing the benefits of self-disclosure, *"I gained a lot from your sharing. I might take further steps to actively disclose my stuttering."*

4.4.5 Speech AI adoption and challenges. Participants report utilizing a range of ASR products for specific use cases in their daily lives. For example, WeChat Voice Messages is commonly used for sending text messages via voice input and converting received voice messages into text. Xiaomi "Xiao Ai" serves purposes such as smart home controls and engaging in casual conversations. iFlytek is primarily used for speech-to-text conversion and daily transcription tasks. Car Navigation Tools enable participants to set destinations and issue navigation commands using voice input. Interestingly, one data collector reflected that some participants reported feeling more comfortable using ASR compared to speaking with real person as they believe ASR would neither judge their speeches nor react differently to their stuttering. Similarly, one participant use ASR to improve fluency and build confidence *"I use that app to practice my speech, such as for the Mandarin proficiency test. On one hand, I do this to desensitize myself, and on the other hand, I feel I need to live up to my job as a teacher."* (P3)

Despite the potential benefits, PWS face unique challenges with ASR products, including recognition errors, time-limited input difficulties, and heightened self-consciousness [21]. Despite these issues, ASR is widely used in China due to its advantages, such as simplifying Chinese typing and improving efficiency. However, usability barriers hinder PWS from leveraging these tools effectively, placing them at a disadvantage in technology use.

5 DISCUSSION

5.1 Technical value of StammerTalk dataset in addressing ASR fluency biases

Representing the disability community adequately and authentically in AI data has been a prominent challenge in AI fairness and accessibility [30, 44]. This challenge is even more pronounced for stuttering, an "invisible" disability that is highly variable, social, and situational. Created by a grassroots stuttering community for AI use, the StammerTalk dataset surpasses existing stuttered speech datasets in its scale, scope, and speech diversity, opening the door for a wide range of technical explorations and interventions for ASR biases.

Although prior research has established stark disparities in ASR model's performance with stuttered versus fluent speech [3, 21, 29], the unprecedented size and diversity of stuttered speech in the StammerTalk dataset will allow deeper understanding of ASR failures across different types of stuttering, stutterers, and speaking contexts. For example, the divergent stuttering patterns captured in unscripted conversation and voice command dictation tasks can inform ASR models about the importance of situational context in understanding stuttered speech, using features exacted from the StammerTalk dataset as a starting point. Also, as previous study on ASR performance with aphasia speech found increased frequency of model hallucination over utterances containing long pauses [17] - a symptom shared by both aphasia and stuttering, targeted analysis on ASR results for different types of stuttering utterances could lead to new insights on common ASR mistakes as well as potential mitigation strategies. Inspecting the types of mistakes made by different ASR models also shed light on the underlying mechanisms within the otherwise blackboxed models that drive their discriminatory behaviors. For example, our results suggest the reliance on language model and semantic context by OpenAI's Whisper model constrains its ability to recognize and transcribe stuttered utterances, while the over-indexing of acoustic features by Meta's wav2vec model could lead to increased homophone errors in its transcript.

Furthermore, the rigorous verbatim transcription, annotated with specific stuttering events, enables ASR systems to recognize and transcribe stuttered utterances as they are, which not only provides a more accurate transcription but also normalizes stuttering in human communications - an attitude clinically proven to benefit people who stutter in the long term [35]. While conventional approach for ASR evaluation routinely remove disfluencies - such as the filler words - from both ground truth and model generated transcript to make it easier to align and compare the reference with

the inference [18], the verbatim transcriptions provided in the StammerTalk dataset allows us to better measure and address fluency biases in ASR models. For example, our audit of the Whisper model using the literal transcript reveals its tendency to artificially “smooth out” stuttered speech in the transcriptions and exposes its embedded ableist biases against speech disfluencies. We also show that such biases can be partially addressed by fine tuning ASR models using the StammerTalk dataset. Compared to the off-the-shelf Whisper model, the fine-tuned model produces more accurate transcriptions of stuttered speech transcriptions with consistent reductions in general and all sub-types of mistakes.

5.2 Social and educational values of collective stuttered voices in Chinese

The StammerTalk dataset also offers unique social and educational values. While speech interfaces and ASR-mediated interactions have been increasingly adopted for convenience, accessibility, and cost-efficiency, the lack of inclusion of users with diverse speech patterns in the research and development of these systems could lead to new accessibility barriers and psychological harms [3, 21, 43]. The StammerTalk dataset can inform HCI researchers about the diversity and variability of speech input, contributing new user personas and design considerations for inclusive speech technologies.

As the only stuttered speech corpus in a non-Western language to our knowledge, the StammerTalk dataset also fills in an important language gap for stuttered speech and opens doors to quantify the linguistic and cultural differences in stuttering between Chinese and other, mostly Eurocentric, languages. Besides, compared to plain text transcripts, the audio format of the personal experiences told by PWS in China creates an intimate channel for self advocacy and empathy building. Listening to personal stories—particularly those highlighting systemic discrimination and psychological struggles- told in stuttered voices —can provide speech AI researchers and designers with a deeper understanding of the goals and needs of PWS, as well as greater awareness of their own fluency biases. On the other hand, for many participants, it was the first time they were able to speak about their stutter and have their stuttered voices heard by the public. As one of the first public discourses about stuttering experiences in China, the dataset provides a platform for collective actions, claiming the much needed space for stuttering in Chinese society.

The dataset also enhances understanding of the social context around stuttering in China, which is essential for creating socially aware and inclusive products. For instance, while products aimed at masking of stuttering have been increasingly rejected by the stuttering community in the US [23, 39, 45], such solutions may appeal to participants in China, where stuttering carries significant personal and professional risks but lacks support infrastructure. Building socially aware products could introduce an ethical dilemma between addressing pressing harms and maintaining fundamental values such as justice and authenticity [37].

5.3 Limitation and future work

Despite the unprecedented scale of the StammerTalk dataset, our work still has several limitations. First, our focus on Chinese stuttered speech restricts its applicability to other languages. Future work could replicate the StammerTalk data collection model across additional languages and dialects, further expanding the diversity and scale of stuttered speech datasets. Additionally, while prior research demonstrates the promise of fine-tuning general ASR models with small amounts of stuttered speech [21, 29], future work could explore how the StammerTalk dataset can advance this direction further. Second, our process focuses on the curation of the dataset but managing the StammerTalk dataset requires significant effort. Although the community intends to open-source their data for scientific and technological advancements, future work should help the community navigate complex legal and technical systems to identify suitable infrastructure for collective ownership, personal data protection, and cross-border data regulations. Finally, while the StammerTalk dataset highlights the trade-offs between technical and social values, future work should explore

strategies to achieve balance. For instance, there may be tensions between selecting topics that resonate deeply with the community and ensuring the diversity of vocabulary required for technical advancements in data collection.

6 CONCLUSION

This work tackles the performance disparity in modern ASR systems for stuttered speech by introducing a community-created, large-scale stuttered speech corpus in Chinese, and demonstrating its effectiveness in benchmarking and diagnosing state-of-the-art ASR models for stuttered speech. Our quantitative and qualitative analysis of the StammerTalk dataset demonstrates the scope and diversity of stuttered utterances it captured, highlighting its unique technical, social and educational value for authentically representing the stuttering community in ASR data.

ACKNOWLEDGMENTS

We would like to extend our heartfelt gratitude to the StammerTalk community, especially Rong Gong, Lezhi Wang, and Jia Bin, for their unwavering support, dedication, and camaraderie throughout this project. Their insights and commitment have been instrumental and inspiring for this work.

REFERENCES

- [1] Sadeen Alharbi, Madina Hasan, Anthony JH Simons, Shelagh Brumfitt, and Phil Green. 2018. A lightly supervised approach to detect stuttering in children’s speech. In *Proceedings of Interspeech 2018*. ISCA, 3433–3437.
- [2] Rosana Ardila, Megan Branson, Kelly Davis, Michael Kohler, Josh Meyer, Michael Henretty, Reuben Morais, Lindsay Saunders, Francis Tyers, and Gregor Weber. 2020. Common Voice: A Massively-Multilingual Speech Corpus. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, Nicoletta Calzolari, Frédéric B chet, Philippe Blache, Khalid Choukri, Christopher Cieri, Thierry Declerck, Sara Goggi, Hitoshi Isahara, Bente Maegaard, Joseph Mariani, H l ne Mazo, Asuncion Moreno, Jan Odijk, and Stelios Piperidis (Eds.). European Language Resources Association, Marseille, France, 4218–4222. <https://aclanthology.org/2020.lrec-1.520>
- [3] Anna Bleakley, Daniel Rough, Abi Roper, Stephen Lindsay, Martin Porcheron, Minha Lee, Stuart Alan Nicholson, Benjamin R Cowan, and Leigh Clark. 2022. Exploring Smart Speaker User Experience for People Who Stammer. In *Proceedings of the 24th International ACM SIGACCESS Conference on Computers and Accessibility*. 1–10.
- [4] O. Bloodstein, N.B. Ratner, and S.B. Brundage. 2021. *A Handbook on Stuttering, Seventh Edition*. Plural Publishing, Incorporated. <https://books.google.com/books?id=Abw0EAAAQBAJ>
- [5] Elaine Blumgart, Yvonne Tran, and Ashley Craig. 2010. Social anxiety disorder in adults who stutter. *Depress. Anxiety* 27, 7 (July 2010), 687–692. <https://doi.org/10.1002/da.20657>
- [6] John Van Borsel, Marie Brepoels, and Janne De Coene. 2011. Stuttering, attractiveness, and romantic relationships: The perception of adolescents and young adults. *Journal of Fluency Disorders* 36, 1 (2011), 41–50.
- [7] Michael P. Boyle. 2018. Enacted stigma and felt stigma experienced by adults who stutter. *Journal of Communication Disorders* 73 (2018), 50–61. <https://doi.org/10.1016/j.jcomdis.2018.03.004>
- [8] Geoffrey A. Coalson, Alexis Crawford, Shanley B. Treleaven, Courtney T. Byrd, Lauren Davis, Lillian Dang, Jillian Ederly, and Alison Turk. 2022. Microaggression and the adult stuttering experience. *Journal of Communication Disorders* 95 (2022), 106180. <https://doi.org/10.1016/j.jcomdis.2021.106180>
- [9] Christopher Constantino, Patrick Campbell, and Sam Simpson. 2022. Stuttering and the social model. *Journal of Communication Disorders* 96 (2022), 106200. <https://doi.org/10.1016/j.jcomdis.2022.106200>
- [10] Christopher Dominick Constantino. 2023. Fostering Positive Stuttering Identities Using Stutter-Affirming Therapy. *Language, Speech, and Hearing Services in Schools* 54, 1 (2023), 42–62. https://doi.org/10.1044/2022_LSHSS-22-00038 arXiv:https://pubs.asha.org/doi/pdf/10.1044/2022_LSHSS-22-00038
- [11] S. Davis, P. Howell, and F. Cooke. 2002. Sociodynamic relationships between children who stutter and their non-stuttering classmates. *Journal of Child Psychology and Psychiatry, and Allied Disciplines* 43, 7 (2002), 939–947.
- [12] Xianghua Ding, Patrick C Shih, and Ning Gu. 2017. Socially embedded work: A study of wheelchair users performing online crowd work in china. In *Proceedings of the 2017 ACM Conference on Computer Supported Cooperative Work and Social Computing*. 642–654.
- [13] Andrews Gavin and MARY HARRIS. 1964. The Syndrome of Stuttering’. *Clinics in developmental Medicine* 17 (1964).
- [14] Jordan R Green, Robert L MacDonald, Pan-Pan Jiang, Julie Cattiau, Rus Heywood, Richard Cave, Katie Seaver, Marilyn A Ladewig, Jimmy Tobin, Michael P Brenner, et al. 2021. Automatic Speech Recognition of Disordered Speech: Personalized Models Outperforming Human Listeners on Short Phrases.. In *Interspeech*. 4778–4782.
- [15] Gerlach H, Totty E, Subramanian A, and Zebrowski P. 2018. Stuttering and Labor Market Outcomes in the United States. *J Speech Lang Hear Res*. 61, 7 (2018), 1649–1663. https://doi.org/10.1044/2018_JSLHR-S-17-0353 <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC6195060/>.
- [16] Peter Howell, Stephen Davis, and Jon Bartrip. 2009. The university college london archive of stuttered speech (uclass). (2009).
- [17] Allison Koenecke, Anna Seo Gyeong Choi, Katelyn X. Mei, Hilke Schellmann, and Mona Sloane. 2024. Careless Whisper: Speech-to-Text Hallucination Harms. In *Proceedings of the 2024 ACM Conference on Fairness, Accountability, and Transparency (Rio de Janeiro, Brazil) (FAccT ’24)*. Association for Computing Machinery, New York, NY, USA, 1672–1681. <https://doi.org/10.1145/3630106.3658996>
- [18] Allison Koenecke, Andrew Nam, Emily Lake, Joe Nudell, Minnie Quartey, Zion Mengesha, Connor Toups, John R Rickford, Dan Jurafsky, and Sharad Goel. 2020. Racial disparities in automated speech recognition. *Proceedings of the national academy of sciences* 117, 14 (2020), 7684–7689.
- [19] Tedd Kourkounakis, Amirhossein Hajavi, and Ali Etemad. 2021. FluentNet: End-to-End Detection of Stuttered Speech Disfluencies With Deep Learning. *IEEE/ACM Trans. Audio, Speech and Lang. Proc.* 29 (sep 2021), 2986–2999. <https://doi.org/10.1109/TASLP.2021.3110146>
- [20] Tedd Kourkounakis, Amirhossein Hajavi, and Ali Etemad. 2021. Fluentnet: End-to-end detection of stuttered speech disfluencies with deep learning. *IEEE/ACM Transactions on Audio, Speech, and Language Processing* 29 (2021), 2986–2999.
- [21] Colin Lea, Zifang Huang, Jaya Narain, Lauren Tooley, Dianna Yee, Dung Tien Tran, Panayiotis Georgiou, Jeffrey P Bigham, and Leah Findlater. 2023. From User Perceptions to Technical Improvement: Enabling People Who Stutter to Better Use Speech Recognition. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems (Hamburg, Germany) (CHI ’23)*. Association for Computing Machinery, New York, NY, USA, Article 361, 16 pages. <https://doi.org/10.1145/3544548.3581224>
- [22] Colin Lea, Vikramjit Mitra, Aparna Joshi, Sachin Kajarekar, and Jeffrey Bigham. 2021. Sep-28k: A Dataset for Stuttering Event Detection from Podcasts with People Who Stutter. <https://arxiv.org/pdf/2102.12394.pdf>
- [23] Jingin Li, Shaomei Wu, and Gilly Leshed. 2024. Re-envisioning Remote Meetings: Co-designing Inclusive and Empowering Videoconferencing with People Who Stutter. In *Proceedings of the 2024 ACM Designing Interactive Systems Conference (Copenhagen, Denmark) (DIS ’24)*. Association for

- Computing Machinery, New York, NY, USA, 1926–1941. <https://doi.org/10.1145/3643834.3661533>
- [24] Qisheng Li and Shaomei Wu. 2024. "I Want to Publicize My Stutter": Community-led Collection and Curation of Chinese Stuttered Speech Data. *Proc. ACM Hum.-Comput. Interact.* 8, CSCW2, Article 475 (Nov. 2024), 27 pages. <https://doi.org/10.1145/3687014>
- [25] Yan Ma, Judith D. Oxley, J. Scott Yaruss, and John A. Tetsnowski. 2023. Stuttering experience of people in China: A cross-cultural perspective. *Journal of Fluency Disorders* 77 (2023), 105994. <https://doi.org/10.1016/j.jfludis.2023.105994>
- [26] Bob MacDonald, Pan-Pan Jiang, Julie Cattiau, Rus Heywood, Richard Cave, Katie Seaver, Marilyn Ladewig, Jimmy Tobin, Michael Brenner, Philip Q Nelson, Jordan R. Green, and Katrin Tomanek. 2021. Disordered Speech Data Collection: Lessons Learned at 1 Million Utterances from Project Euphonia.
- [27] Valentin Mendelev, Tina Raissi, Guglielmo Camporese, and Manuel Giollo. 2021. Improved robustness to disfluencies in rnn-transducer based speech recognition. In *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 6878–6882.
- [28] Meta. 2023. Speech Fairness Dataset. <https://ai.meta.com/datasets/speech-fairness-dataset/>.
- [29] Vikramjit Mitra, Zifang Huang, Colin Lea, Lauren Tooley, Panayiotis Georgiou, Sachin Kajarekar, and Jefferey Bigham. 2021. Analysis and Tuning of a Voice Assistant System for Dysfluent Speech. <https://arxiv.org/pdf/2106.11759.pdf>
- [30] Joon Sung Park, Danielle Bragg, Ece Kamar, and Meredith Ringel Morris. 2021. Designing an Online Infrastructure for Collecting AI Data From People With Disabilities. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency (Virtual Event, Canada) (FAcT '21)*. Association for Computing Machinery, New York, NY, USA, 52–63. <https://doi.org/10.1145/3442188.3445870>
- [31] Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. 2023. Robust speech recognition via large-scale weak supervision. In *Proceedings of the 40th International Conference on Machine Learning (Honolulu, Hawaii, USA) (ICML '23)*. JMLR.org, Article 1182, 27 pages.
- [32] Nan Bernstein Ratner and Brian MacWhinney. 2018. Fluency Bank: A new resource for fluency research and practice. *Journal of fluency disorders* 56 (2018), 69–80.
- [33] Johnny Saldaña. 2021. *The coding manual for qualitative researchers*. sage.
- [34] Olabanji Shonibare, Xiaosu Tong, and Venkatesh Ravichandran. 2022. Enhancing ASR for Stuttered Speech with Limited Data Using Detect and Pass. arXiv:2202.05396 [eess.AS]
- [35] Vivian Sisskin. 2023. Disfluency-Affirming Therapy for Young People Who Stutter: Unpacking Ableism in the Therapy Room. *Language, Speech, and Hearing Services in Schools* 54, 1 (2023), 114–119.
- [36] Joshua St. Pierre. 2012. The Construction of the Disabled Speaker: Locating Stuttering in Disability Studies. *Canadian Journal of Disability Studies* 1, 3 (Aug. 2012), 1–21. <https://doi.org/10.15353/cjds.v1i3.54>
- [37] Sharifa Sultana, François Guimbretière, Phoebe Sengers, and Nicola Dell. 2018. Design Within a Patriarchal Society: Opportunities and Challenges in Designing for Rural Women in Bangladesh. In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems (Montreal QC, Canada) (CHI '18)*. Association for Computing Machinery, New York, NY, USA, 1–13. <https://doi.org/10.1145/3173574.3174110>
- [38] Seth E Tichenor and J Scott Yaruss. 2021. Variability of stuttering: Behavior and impact. *American Journal of Speech-Language Pathology* 30, 1 (2021), 75–88.
- [39] Space to Stutter. 2024. Open Letter to Samsung. <https://www.spacetostutter.org/impulse>. Accessed on 1/23/2025..
- [40] Jimmy Tobin and Katrin Tomanek. 2022. Personalized automatic speech recognition trained on small disordered speech datasets. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 6637–6641.
- [41] Seda Türkili, Serkan Türkili, and Zeynep Feryal Aydın. 2022. Mental well-being and related factors in individuals with stuttering. *Heliyon* 8, 9 (Sept. 2022), e10446. <https://doi.org/10.1016/j.heliyon.2022.e10446>
- [42] UNIVERSITY OF ILLINOIS URBANA-CHAMPAIGN. 2023. Speech Accessibility Project. <https://speechaccessibilityproject.beckman.illinois.edu/>. Accessed on 1/21/2024..
- [43] Kimi Wenzel, Nitya Devireddy, Cam Davison, and Geoff Kaufman. 2023. Can Voice Assistants Be Microaggressors? Cross-Race Psychological Responses to Failures of Automatic Speech Recognition. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems (CHI '23)*. Association for Computing Machinery, New York, NY, USA, Article 109, 14 pages. <https://doi.org/10.1145/3544548.3581357>
- [44] Meredith Whittaker, Cynthia L. Bennett Meryl Alper, Sara Hendren, Elizabeth Kazianas, Mara Mills, Meredith Ringel Morris, Joy Lisi Rankin, Emily Rogers, Marcel Salas, and Sarah Myers West. 2019. Disability, Bias & AI Report. *AI Now Institute* (20 11 2019).
- [45] Shaomei Wu. 2023. "The World is Designed for Fluent People": Benefits and Challenges of Videoconferencing Technologies for People Who Stutter. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*. 1–17.
- [46] yeyupiaoling. [n. d.]. Whisper-Finetune. <https://github.com/yeyupiaoling/Whisper-Finetune>. Accessed: 2025-01-22.
- [47] Megan M Young, Courtney T Byrd, Rodney Gabel, and Andrew Z White. 2022. Self-disclosure experiences of adults who stutter: An interpretative phenomenological analysis. *Am. J. Speech. Lang. Pathol.* 31, 5 (Sept. 2022), 2045–2060. https://doi.org/10.1044/2022_AJSLP-22-00048
- [48] Xin Zhang, Iván Vallés-Pérez, Andreas Stolcke, Chengzhu Yu, Jasha Droppo, Olabanji Shonibare, Roberto Barra-Chicote, and Venkatesh Ravichandran. 2022. Stutter-TTS: Controlled synthesis and improved recognition of stuttered speech. *arXiv preprint arXiv:2211.09731* (2022).

A ADDITIONAL DESCRIPTIVE STATISTICS OF STAMMERTALK DATASET

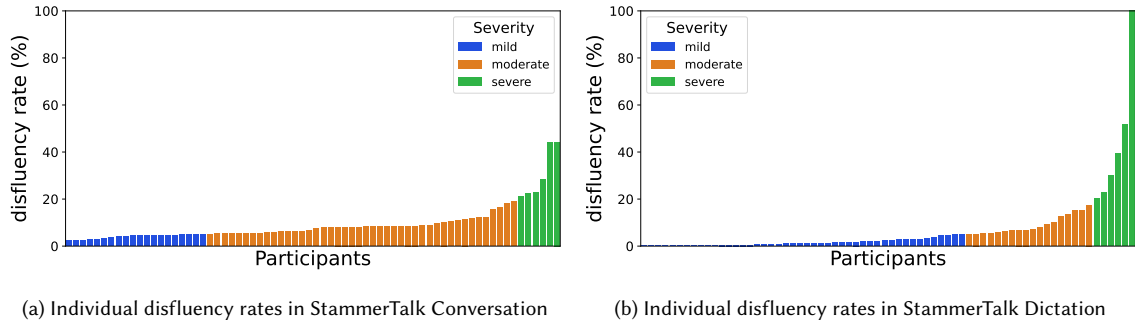


Fig. 4. Disfluency rates of all 70 participants in Conversation and Dictation tasks, sorted from low to high, and categorized into mild (0-5%), moderate (5-20%), and severe (20%+) stuttering groups.

B ADDITIONAL BENCHMARKING RESULTS

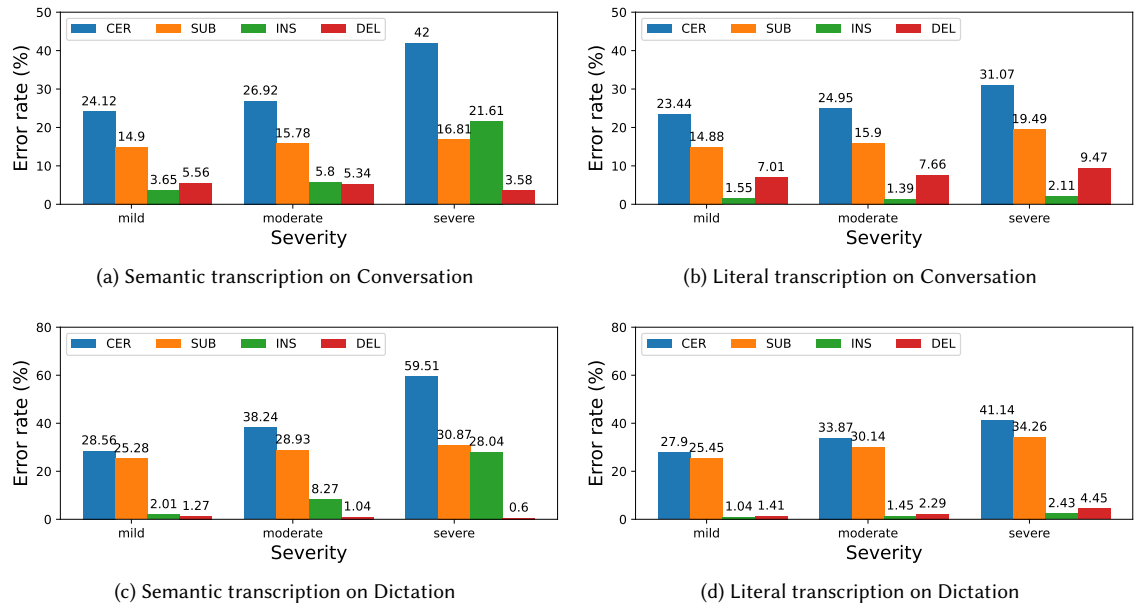


Fig. 5. Character error rate (CER), substitution (SUB), insertion (INS) and deletion (DEL) error rates for fine-tuned Wav2Vec 2.0 evaluated on StammerTalk Conversation and Dictation

Table 3. Whisper model "smooths" the transcriptions by removing words with low semantic value, as indicated by the underlined characters.

Reference	Whisper model output
就在那个继续深造的也有 嗯我觉得深圳 <u>他</u> 到处 <u>他</u> 都是花钱的地方就是吃喝玩乐 <u>他</u> 肯定是	就在继续深造的也有 我觉得深圳到处都是花田的地方吃喝玩乐肯定是

Table 4. Three examples of utterances from a severe PWS, characterized by frequent word repetitions. The wav2Vec model produced homophone substitutions, as indicated by the underlined characters.

Annotation	Wav2Vec model output
当[<u>当</u> 当]时我上/b[<u>上</u>]去的时候	当档单舍瓦上上去的时候
我/b[<u>现</u> 现]就[<u>就</u>]挺自/r[<u>卑</u>]的	我先线先千就就点自杯给的
呃/i[<u>进</u>]行那个自[<u>自</u>]我介[<u>介</u>]绍，呃/i	而 <u>仅仅</u> 进行了个自自我 <u>界</u> 介绍和儿

C CODING SCHEME

Category	Subcategory	Codes
1. Stuttering Experiences	1.1 Personal Experiences	Fear, anxiety; Tension, feeling out of control; Emotional responses after stuttering (e.g., embarrassment, frustration); Mental health problems (e.g., depression, anxiety disorder)
	1.2 Social Experiences	Negative reactions from others (e.g., mockery, making fun of, distrust); Discrimination or prejudice; Positive experience: Supportive social interactions
	1.3 Misconceptions of Stuttering	Stuttering is from imitation; Stuttering can be cured; Stuttering is because of nervousness; "If you speak slowly, you wouldn't stutter"
	1.4 Attitudes Toward Stuttering	Acceptance; Not accepting and wanting to cure
2. Coping Mechanisms	2.1 Strategies for Managing Stuttering	Avoidance behaviors (e.g., avoiding certain words, speaking situations, professions, relationships); Speech therapy techniques (e.g., fluency shaping, deep breath); Speech therapy program; Existing resources (e.g., books, online); Seeking help from SLP
	2.2 Emotional Coping	Internal dialogues (e.g., self-reassurance, self-acceptance); Seeking support from family and friends; Seeking support from online communities; Seeking support from mental health professionals
3. Self-Disclosure	3.1 Levels of Disclosure	Public disclosure; Disclosure to family and close friends; Disclosure to co-workers or in professional settings; No explicit self-disclosure (assume others are aware)
	3.2 Barriers to Self-Disclosure	Fear of stigma; Previous negative experiences (e.g., dismissal, denial of stuttering)
	3.3 Challenges After Disclosure	Misguided advice; Negative reactions (e.g., mockery, denial); Social and professional consequences
4. Occupational Distribution and Challenges	4.1 Occupation	Coding by type (e.g., teacher, doctor, customer service)
	4.2 Work-Related Challenges	Communication-intensive roles (e.g., teaching, public speaking); Impact on career progression or job opportunities
	4.3 Urban vs. Cosmopolitan vs. Overseas Experiences	Coding by area if any

Category	Subcategory	Codes
5. ASR Products and Usage	5.1 Product Usage	Frequency of use (e.g., daily, occasionally); Frequent use at home or in solo settings; Avoidance in public or group settings; Purpose: Speech-to-text, Control smart home devices, Accent reduction, Fluency shaping
	5.2 Challenges	Recognition issues: inaccurate recognition of the stuttered speech; Premature cancellation of input if there are delays; Error-prone results requiring manual corrections
	5.3 Social Dynamics	Hesitation or avoidance of ASR in front of others; Preference for using ASR in private
6. Dynamic Nature of Stuttering	6.1 Variability Over Time	Changes in stuttering severity across life stages; Impact of specific situations (e.g., stress, public speaking)
	6.2 Contextual Factors	Variations in stuttering based on the audience or setting; External triggers or mitigators (e.g., pressure, comfort levels)