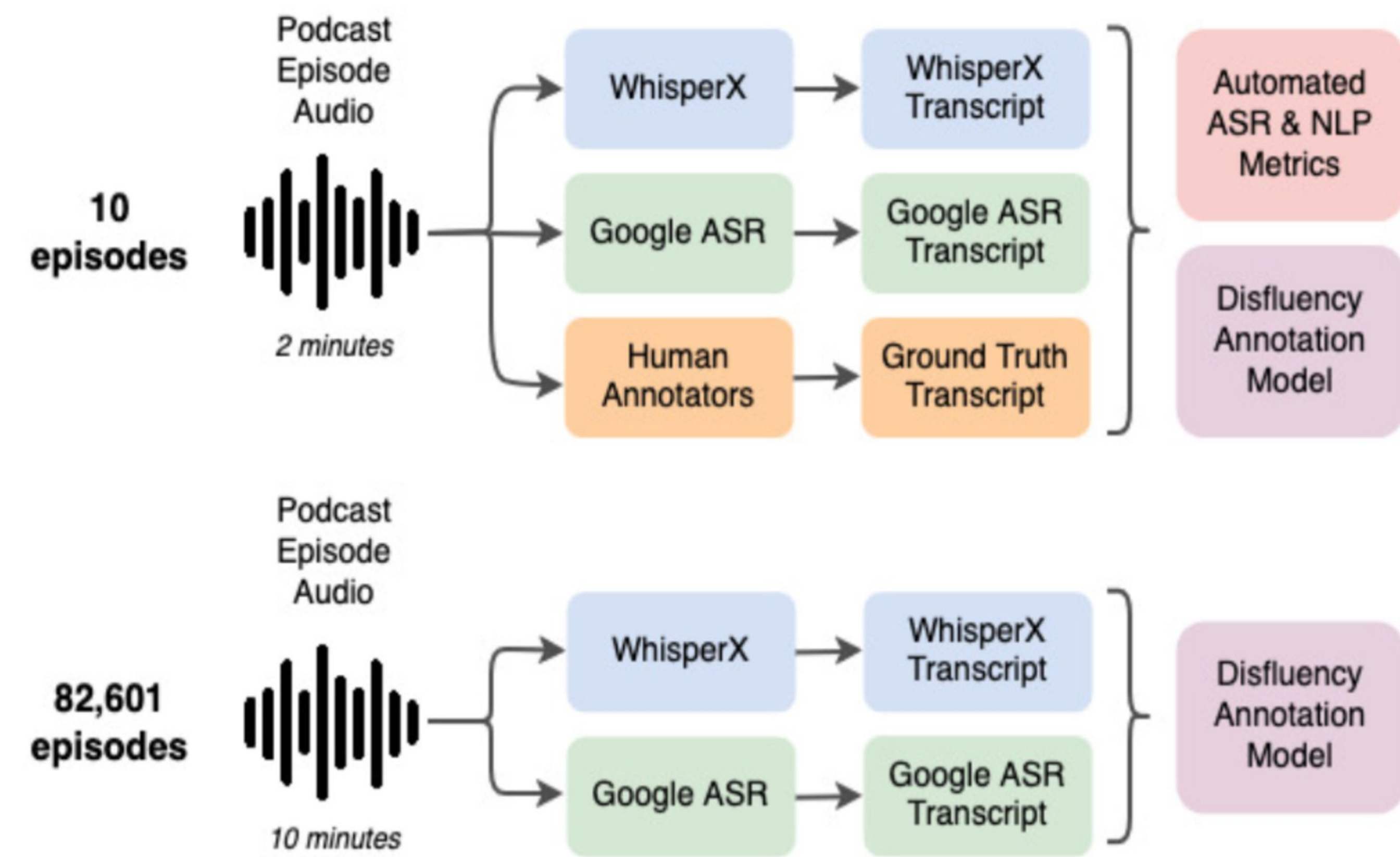


# Comparing ASR Systems in the Context of Speech Disfluencies

Maria Teleki, Xiangjue Dong, Soohwan Kim, James Caverlee  
Texas A&M University

## Introduction & Experimental Settings

- We evaluate the disfluency capabilities of two ASR systems – **WhisperX** [1] and **Google ASR** [2] – in terms of their interactions with a **parsing-based disfluency annotation model** [3].
  - Why? It's natural for application developers to plug an ASR-created transcript into a disfluency removal model.
  - 4-6% of non-scripted speech is disfluent [4].
  - We use [3] to annotate the **3 types of disfluencies**:
    - interjections (**INTJ**) – ex: *let's go to the uh store today*
    - parentheticals (**PRN**) – ex: *let's to go the store, wait no, the movies today*
    - edited nodes (**EDITED**) – ex: *let's to go the store, wait no, the movies today*
- We use the **Spotify Podcasts Dataset** [5] for our analysis.
  - We obtain **ground truth transcripts via human annotations (N=3)**.



### RQ1 (10 episodes): How does the choice of ASR system (WhisperX, Google ASR) impact performance (as compared to the human-annotated ground truth)?

- While WhisperX performs better overall in terms of automated metrics, Google ASR outperforms WhisperX in WIL and BLEU for specifically non-scripted podcasts.

		Character-level	Word-level		Sentence-level		
		CER (↓)	WER (↓)	WIL (↓)	ROUGE-L (↑)	BERTScore (↑)	BLEU (↑)
Scripted	Google ASR	3.46±2.07	7.39±2.99	15.02±1.67	93.83±2.46	97.66±1.31	85.09±0.32
	WhisperX	1.87±1.49	3.36±1.37	14.01±2.45	97.41±0.93	99.03±0.53	86.24±1.12
Non-Scripted	Google ASR	8.87±5.95	12.98±6.96	15.03±0.67	90.48±5.06	96.29±2.07	84.85±1.56
	WhisperX	6.05±3.77	9.74±5.32	15.32±0.97	93.34±3.29	97.40±1.25	84.71±1.96
All	Google ASR	6.71±5.37	10.47±6.18	15.02±1.09	91.82±4.39	96.84±1.86	84.95±1.18
	WhisperX	4.38±3.64	7.19±5.21	14.79±1.73	94.97±3.27	98.05±1.30	85.32±1.78

### RQ2 (10 episodes): How does the choice of ASR system (WhisperX, Google ASR) impact the specific disfluency types which are transcribed (as compared to the human-annotated ground truth)?

- WhisperX transcribes closer to the ground truth number of *uhs*, *ums*, and INTJ nodes than Google ASR.
  - The ground truth number of *uhs* and *ums* is higher.
- Google ASR transcribes closer to the ground truth number of EDITED nodes than WhisperX.
- WhisperX and Google ASR transcribe the same number of PRN nodes.

		C <sub>“uh”</sub>	C <sub>“um”</sub>	C <sub>INTJ</sub>	C <sub>PRN</sub>	C <sub>EDITED</sub>
		Scripted	Ground Truth	0	0	1.00±0.82
Scripted	Google ASR	0	0	1.00±0.82	0	1.50±1.91
	WhisperX	0	0	0.75±0.96	0	0.75±1.50
Non-Scripted	Ground Truth	1.67±1.97	1.33±1.21	9.06±6.81	2.00±2.38	5.33±4.25
	Google ASR	0	0	6.33±5.32	2.17±2.93	5.33±2.50
Non-Scripted	WhisperX	0.33±0.82	0.67±0.82	7.83±6.40	2.17±2.40	3.67±2.73
	Ground Truth	1.00±1.70	0.80±1.14	5.83±6.59	1.30±2.02	3.43±4.04
All	Google ASR	0	0	4.20±4.85	1.30±2.45	3.80±2.94
	WhisperX	0.20±0.63	0.40±0.70	5.00±6.04	1.30±2.11	2.50±2.68

### RQ3 (82,601 episodes): Are these findings consistent at a large scale?

- Google ASR transcribes hardly any *uhs* or *ums*.
- Same trend as small-scale in RQ2: WhisperX transcribes more INTJ nodes, while Google ASR transcribes more EDITED nodes – however Google ASR transcribes more PRN nodes.
  - We hypothesize this is due to the *vocabulary diversity of WhisperX versus that of Google ASR*.

	C <sub>“uh”</sub>	C <sub>“um”</sub>	C <sub>INTJ</sub>	C <sub>PRN</sub>	C <sub>EDITED</sub>
Google ASR	0.09±0.35	0.25±0.70	48.02±37.12	12.71±11.29	30.26±13.71
WhisperX	1.38±3.03	1.69±3.14	50.90±39.88	10.84±9.79	16.71±9.58

## References

- M. Bain, et al., “WhisperX: Time Accurate Speech Transcription of Long-Form Audio,” in Interspeech, 2023.
- Google Cloud, “Speech-To-Text: Automatic Speech Recognition,” 2024. [Online]. Available: <https://cloud.google.com/speech-to-text>
- P. J. Lou and M. Johnson, “Improving Disfluency Detection by Self-Training a Self-Attentive Model,” in ACL, 2020.
- E. Shriberg, “Preliminaries to a Theory of Speech Disfluencies,” Ph.D. dissertation, 1994.
- A. Clifton, et al., “100,000 Podcasts: A Spoken English Document Corpus,” in COLING, 2020.

## Conclusion

These results suggest that it may be **beneficial to select an ASR system based on the distribution of disfluent node types present in the data.**

