

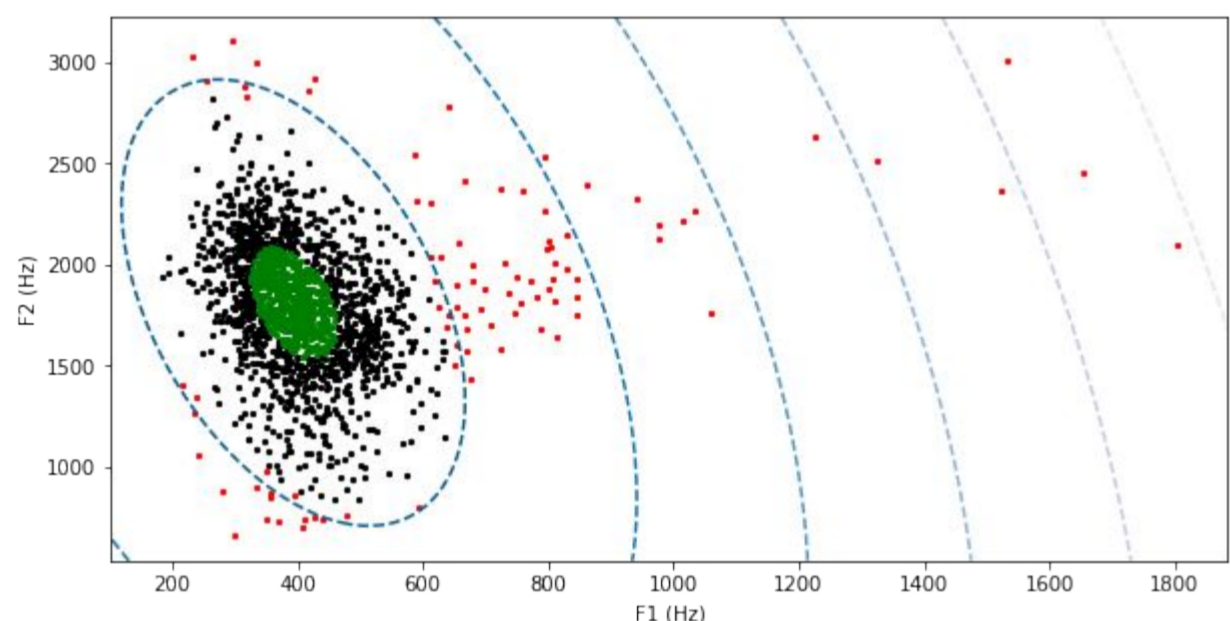
Emily P. Ahn, Gina-Anne Levow, Richard A. Wright, & Eleanor Chodroff
 University of Washington, USA University of Zurich, CH

1. Motivation

- > Automatic forced alignment and phonetic measurement aids field linguists, phoneticians, sociolinguists
- > Understand outliers from a fully automated corpus phonetics pipeline
 - Distinguish between technical errors & true linguistic variation
 - Develop taxonomy of error types

2. Methods

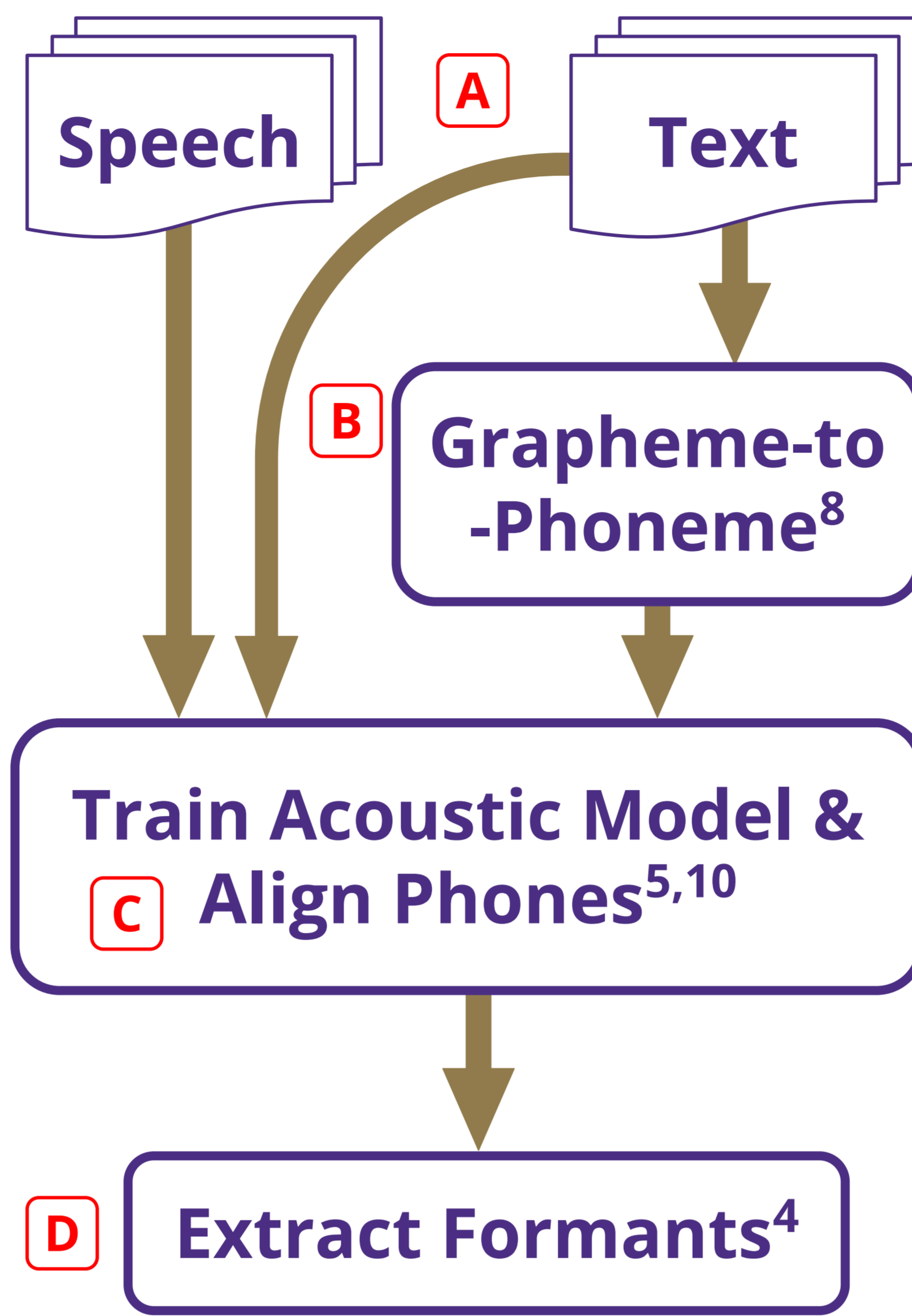
- Download 2 read speech corpora**
 - Wilderness³ Bible & VoxClamantis⁹ alignments/formants
 - Mozilla Common Voice² sentences & VoxCommunis¹ alignments/formants
- Discover vowel formant outliers**
 - Mahalanobis distance



- Annotate with new taxonomy**
 - 840 vowel samples (600 outliers, 240 near-means)
 - 5 trained linguists (Krippendorff's alpha = 0.86, strong agreement)

	Available Corpus		Analyzed Corpus						# Outliers	% Outliers
	Total Hours	Total Spkrs	Hours	Spkrs	Vowel Utts	Vowel Types	Vowel Tokens			
Wilderness										
Hausa	20:40	5+	20:40	5+	9626	5	303577	9698	3.19%	
Kazakh	18:51	5+	18:51	5+	8085	6	204701	22148	10.82%	
Swedish	16:46	1	16:46	1	9516	16	204701	15106	8.28%	
Common Voice v8										
Hausa	3:23	17	0:57	8	772	5	11490	583	5.07%	
Kazakh	1:27	72	1:06	46	796	11	10967	642	5.85%	
Swedish	39:28	674	1:02	203*	1000*	16	11230	513	4.57%	

3a. Pipeline



3b. Assumptions

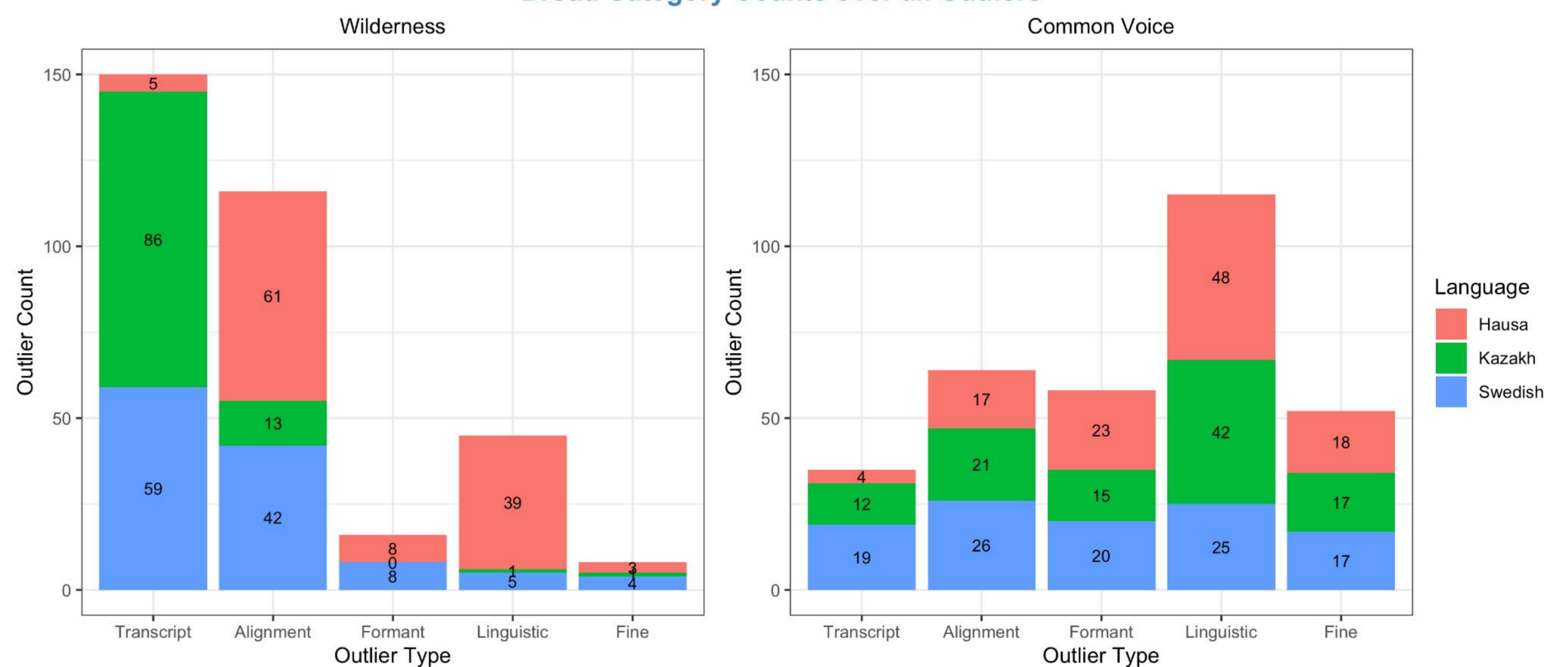
- Script = Transcript
- Phonetic transcription is accurate
- Segmentation is accurate
- Acoustic-phonetic measurement is accurate

3c. Taxonomy

- Transcript Error**
 - Extra sounds (phones, syllables)
 - Extra transcript
 - Broad mismatch
- Alignment Error**
 - Target overlap
 - Broad alignment issue
- Formant Error**
 - Tracker and formant Hz mismatch
- Linguistic Variation**
 - Deletion of target vowel
 - Change (different vowel produced)
- Fine**

4. Results

Broad Category Counts over all Outliers



5. Case Studies

Why are these vowel formants outlying?

→ G2P specification & output may not explicitly represent pronunciation

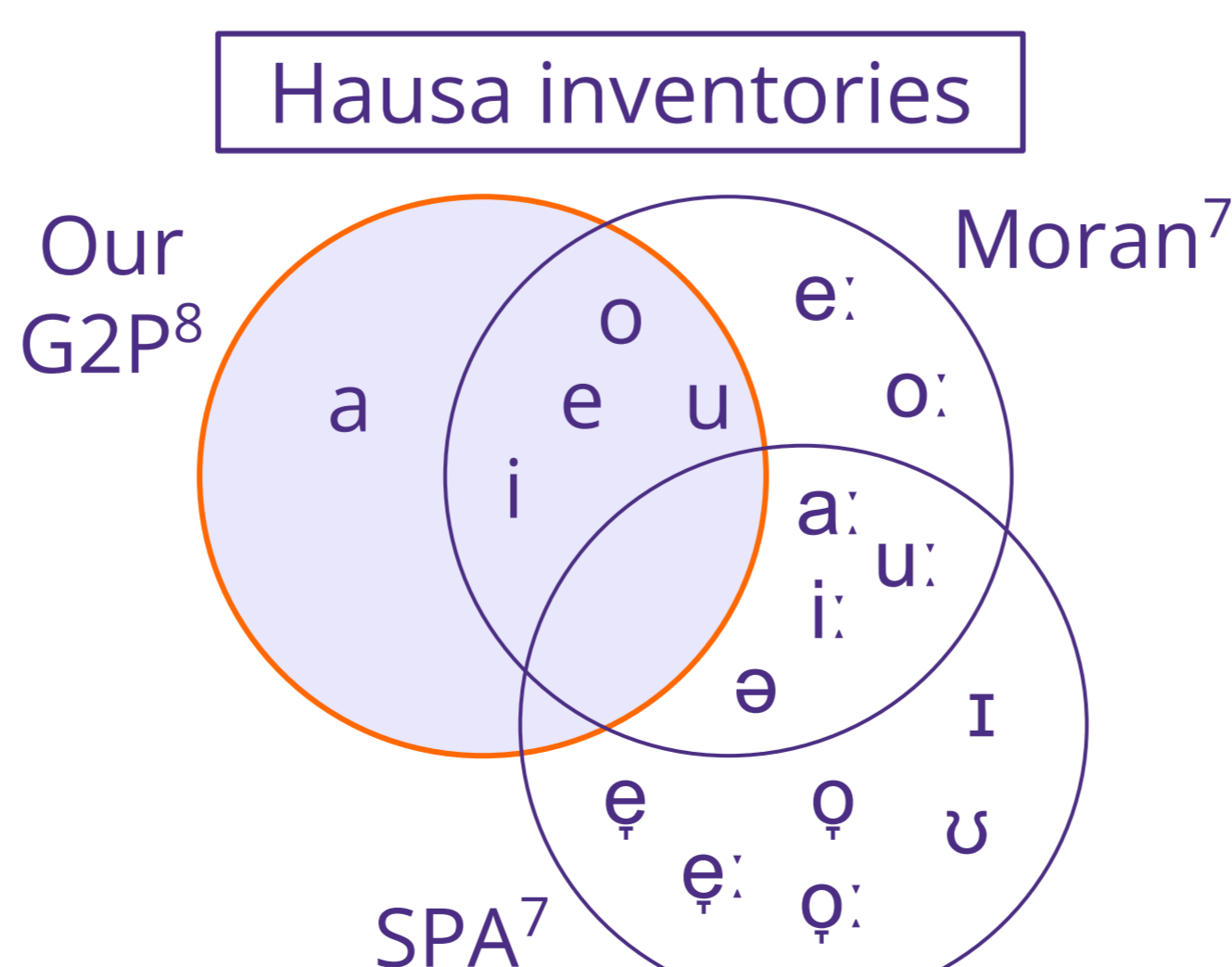
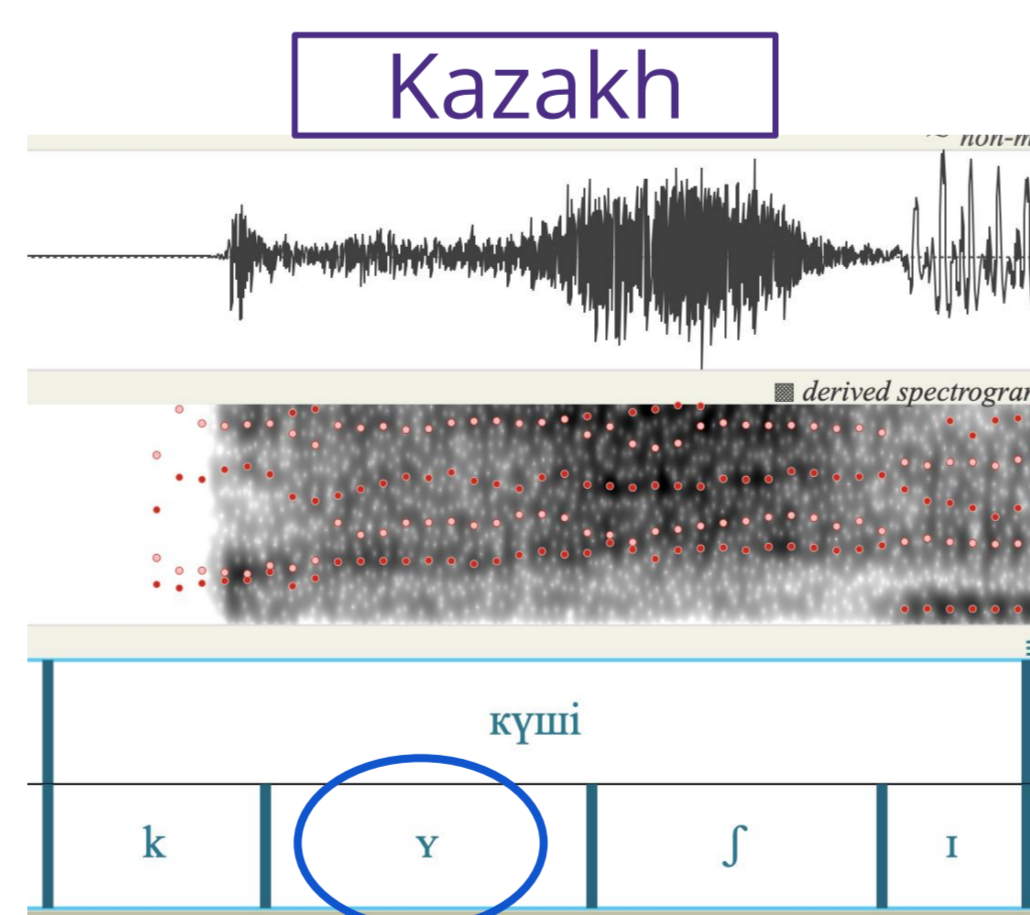
1) High Vowel Deletion in Kazakh

High (short) vowels in Kazakh are more susceptible to reduction (McCollum & Chen, 2021).

→ High vowels are 1.7 times more likely to be deleted than non-high vowels ($p < 0.001$), especially in Kazakh ($p < 0.001$)

2) Vowel Length in Hausa

- > 44% of Hausa outliers & 64% of Hausa near-means marked as Linguistic Change
- > Issue: disagreement of phoneme inventory
 - o In PHOIBLE⁷, some linguists include long vowels, while our G2P does not
 - Vowel length semi-predictable, phonemic?
 - Vowel quality more centralized?



6. Conclusion

Summary

- We develop a novel outlier taxonomy for a corpus phonetics pipeline
- The distribution of outliers reveals dataset quality & quirks, language-specific phenomena

Future Work

- > Apply taxonomy to new data, measures
- > Identify efficient solutions to avoid & correct errors

References

- Ahn, E., & Chodroff, E. (2022). VoxCommunis: A corpus for cross-linguistic phonetic analysis. In *LREC*.
- Ardila et al. (2020). Common Voice: A massively-multilingual speech corpus. In *LREC*.
- Black, A. W. (2019). CMU wilderness multilingual speech dataset. In *ICASSP*.
- Boersma, B. & Weenink, D. (2019). Praat: Doing Phonetics by Computer (Version 6.0.16).
- McAuliffe et al. (2017). Montreal Forced Aligner: Trainable text-speech alignment using Kaldi. In *Interspeech*.
- McCollum, A. G. and Chen, S. (2021). Kazakh. In *JIPA*.
- Moran, S. & McCloy, D., Eds. (2019). PHOIBLE 2.0.
- Mortensen et al. (2018). EpiTran: Precision G2P for many languages. In *LREC*.
- Salesky et al. (2020). A corpus for large-scale phonetic typology. In *ACL*.
- Wiesner et al. (2019). Zero-shot pronunciation lexicons for cross-language acoustic model transfer. In *IEEE ASRU*.

Acknowledgments

We thank Anna Batra, Sam Briggs, Ivy Guo, and Emma Miller. This work was supported in part by NSF GRFP grant DGE-2140004 and SNSF grant 208460.