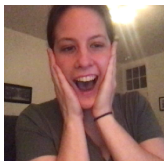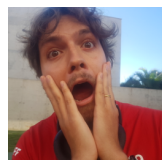# A Corpus For Large-Scale Phonetic Typology
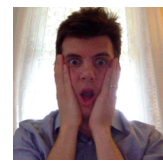
Elizabeth Salesky    Eleanor Chodroff    Tiago Pimentel    Matthew Wiesner
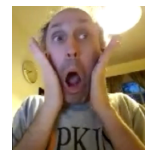
VoxClamantis in deserto:
"a voice crying out in the wilderness"

Ryan Cotterell        Alan W Black        Jason Eisner

1

'in the beginning'
**English**
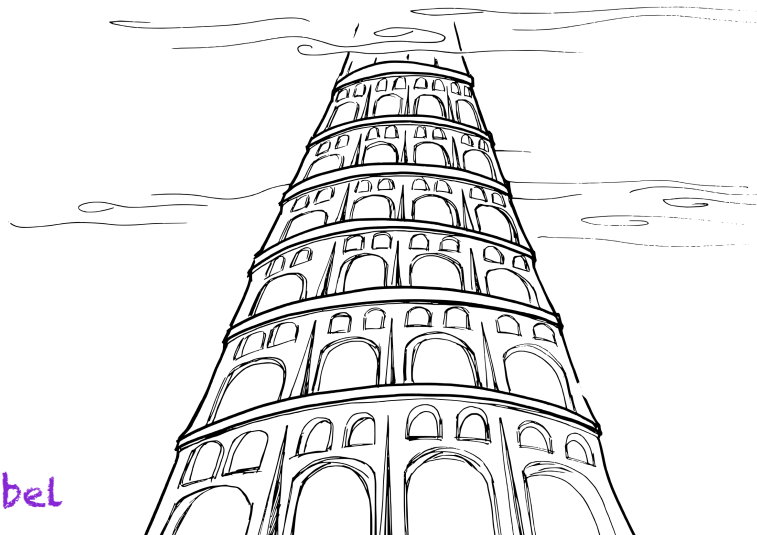
'ipeuhcan'
**Nahuatl**

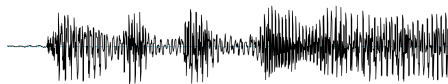'am Anfang'
**German**

'በመጀመሪያ'
**Amharic**

In the beginning, there was **SPEECH**

Tower of Babel

'in the beginning'
**English**

'ipeuhcan'
**Nahuatl**

'am Anfang'
**German**

'በመጀመሪያ'
**Amharic**

# In the beginning, there was **SPEECH**

## Then the linguist asked:

How do speech and language vary?

↪ **prior cross-linguistic phonetic studies have relied on reported [language-aggregate] measurements**

## We create our new corpus, VoxClamantis v1.0, to answer this question!

✔ spoken readings of the Bible
✔ >600 languages
✔ time-aligned phonemic transcriptions
✔ phonetic measures for vowel and sibilant **tokens**

# This talk

① **WHY** we want this data

② **HOW** we create it

③ **CASE STUDIES** validating the corpus & illustrating two possible uses

# Why?

# Variation in and across languages

⑤
**Spanish**

⑦
**Romanian**

| Spanish | Romanian |
|---------|----------|
| /i/ | /i/ |
| /u/ | /u/ |
| /o/ | /o/ |
| /e/ | /e/ |
| /a/ | /a/ |
|  | /ɨ/ |
|  | /ə/ |

We know phonetic variation within a language, but what are its range and limits?
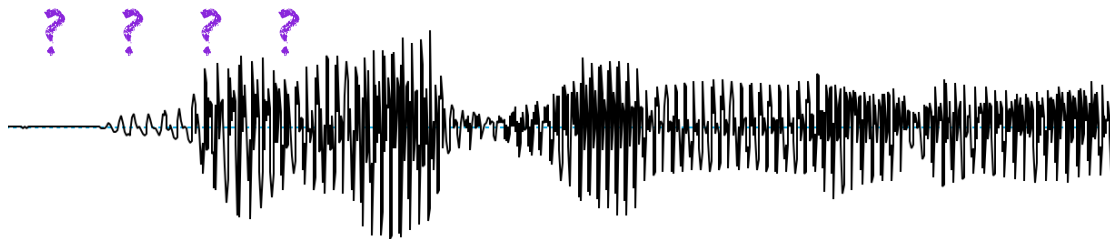
How does the number and set of phonemic categories influence their realizations?

6

# How?

① speech
② transcripts
③ phonemic labels

**Amharic**

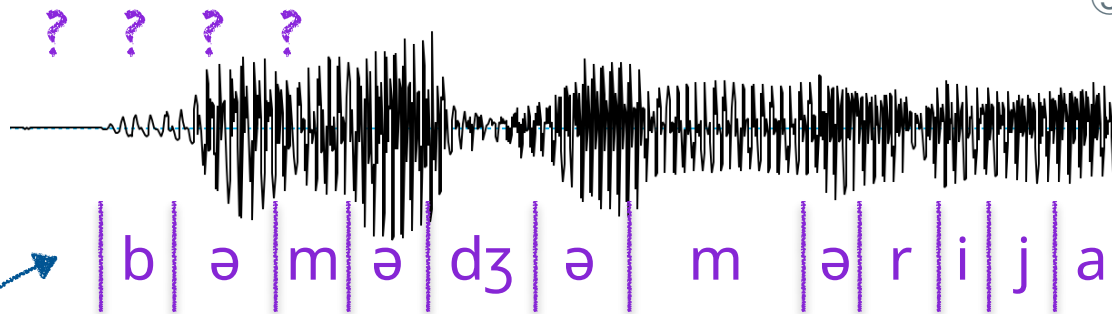? ? ? ?

bəmədʒəmərija

በመጃመሪያ

Grapheme-to-Phoneme

(G2P)

① speech
② transcripts
③ phonemic labels
④ time alignments
⑤ phonetic measures

**Amharic**

? ? ? ?

| b | ə | m | ə | dʒ | ə | m | ə | r | i | j | a

Forced alignment

(HMM acoustic model)

በመጃመሪያ

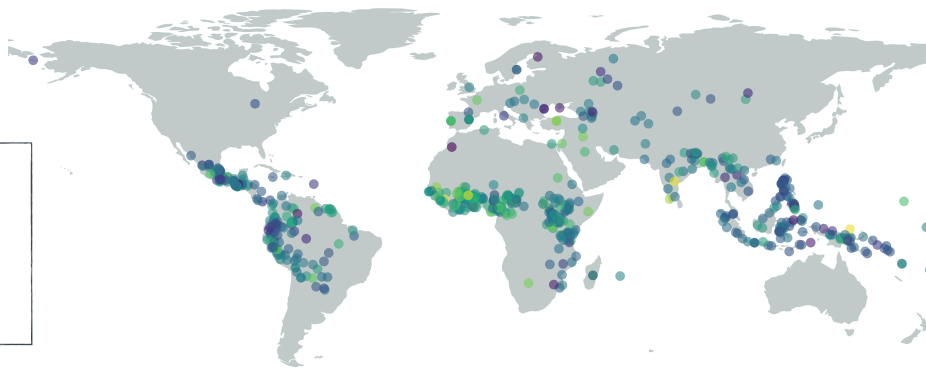Phonetic measures (R or Praat):

Formant frequencies, mid-frequency peak, duration...

① speech
② transcripts

CMU Wilderness
(2019)

699 Bible readings!

with ① speech!

and ② transcripts!
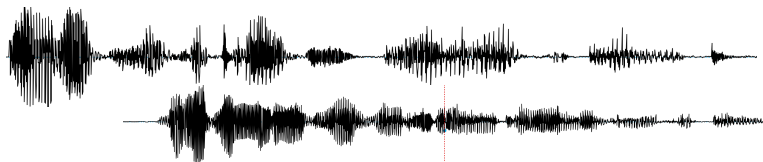
'በመጃመሪያ'
Amharic

>1TB 😱

>6 years of CPU compute 😱

① speech
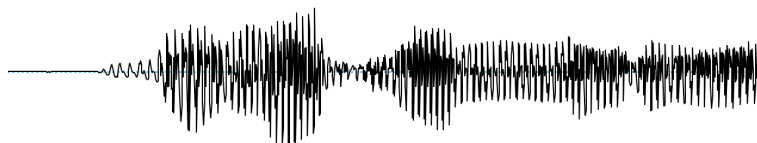② transcripts
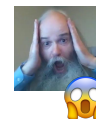
# CMU Wilderness dataset

**Chapter:**

~30min

**1 የፍጥረት አጀማመር** በመጀመሪያ እግዚአብሔር (ኤሎሃም) ሰማያትንና ምድርን ፈጠረ። 2 ምድርም ቅርጽ የለሽና ባዶ ነበረች።※ የምድርን ጥልቅ ስፍራ ሁሉ ጨለማ ውጦት ነበር። የእግዚአብሔርም (ኤሎሃም) መንፈስ በውሆች ላይ ደረብብ ነበር። 3 ከዚያም እግዚአብሔር (ኤሎሃም) "ብርሃን ይሁን" አለ፤ ብርሃንም ሆነ። 4 እግዚአብሔርም (ኤሎሃም) ብርሃኑ መልካም እንደሆነ አየ፤ ብርሃኑን ከጨለማ ለየ። 5 እግዚአብሔርም (ኤሎሃም) ብርሃኑን "ቀን"፣ ጨለማውን "ሌሊት" ብሎ ጠራው። መሽ፤ ነጋም፤ የመጀመሪያ ቀን። 6 እግዚአብሔር (ኤሎሃም)፣ "ውሃን ከውሃ የሚለይ ጠፈር በውሆች መካከል ይሁን" አለ። 7 ስለዚህ እግዚአብሔር (ኤሎሃም) ጠፈርን አድርጎ ከጠፈሩ በላይና ከጠፈሩ በታች ያለውን ውሃ ለየ፤ እንዲለውም ሆነ። 8 እግዚአብሔር (ኤሎሃም) ጠፈርን "ሰማይ" ብሎ ጠራው። መሽ፤ ነጋም፤ ሁለተኛ ቀን። 9 ከዚያም እግዚአብሔር (ኤሎሃም)፣ "ከሰማይ በታች ያለው ውሃ በአንድ.
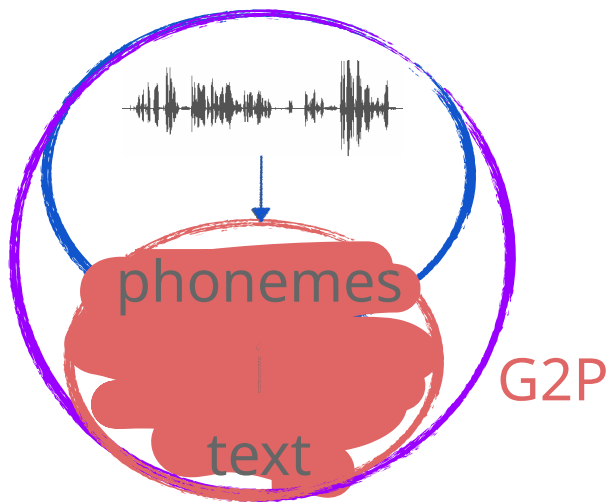
. . .

**Utterance:**

<30s 😱

በመጀመሪያ

11

# Phoneme "Transcriptions"— Grapheme-to-Phoneme

**① Linguist-created rules** *(Epitran)*

**39 readings**

690  64

**② Wisdom of Crowds** *(Wiktionary/WikiPron)*

*+ our own WFST-models (Phonetisaurus* 🦕 *)*

**18 readings** *(disjoint)*

690  165

**③ Naïve baseline** *(Unitran)*

😱 "first-pass transcription"

**All 690 readings**

690

57 readings
**"High-resource (HR)"**

ALL 690 readings
**"First-pass (FP)"**

🤔 why provide **FP** alignments for languages with **HR** ? We'll come back to that 😉

**Amharic**

**bəmədʒəmərija**

Forced alignment

(HMM acoustic model)

① speech
② transcripts
③ phonemic labels
④ time alignments

Amharic

? ? ? ?

|b|ə|m|ə|dʒ|ə| m |ə|r|i|j|a
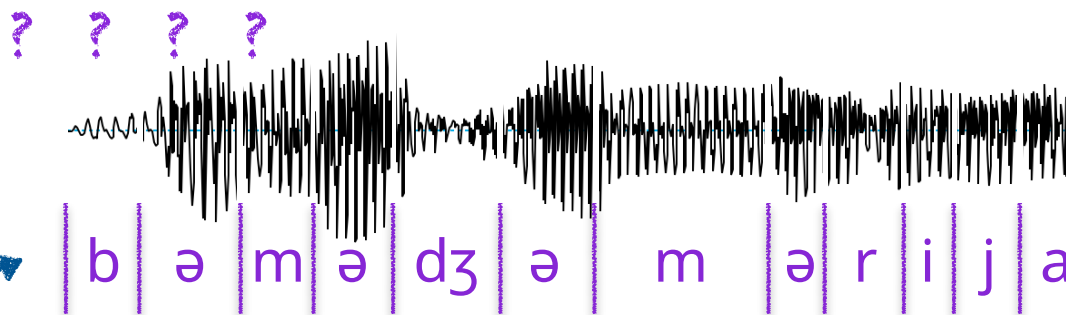
Forced alignment

(HMM acoustic model)

|b|

start     end
time     time

16

① speech
② transcripts
③ phonemic labels
④ time alignments

**Amharic**

? ? ? ?

| b | ə | m | ə | dʒ | ə | m | ə | r | i | j | a |

Forced alignment

(HMM acoustic model)

| b |

start time    end time

17

① speech
② transcripts
③ phonemic labels
④ time alignments

**Amharic**

**Phoneme tokens:**

b

ə
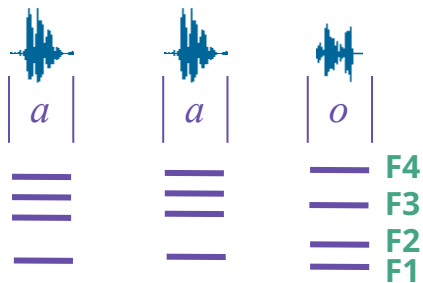
m

...

b

start time    end time

18

# Phonetic Measures

① speech
② transcripts
③ phonemic labels
④ time alignments
⑤ phonetic measures

## VOWELS

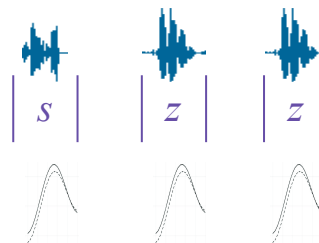## SIBILANTS

$a$ | $a$ | $o$

F4
F3
F2
F1

$s$ | $z$ | $z$

*eg high-amplitude frequencies*

**Formants**

**PRAAT TEXTGRID**

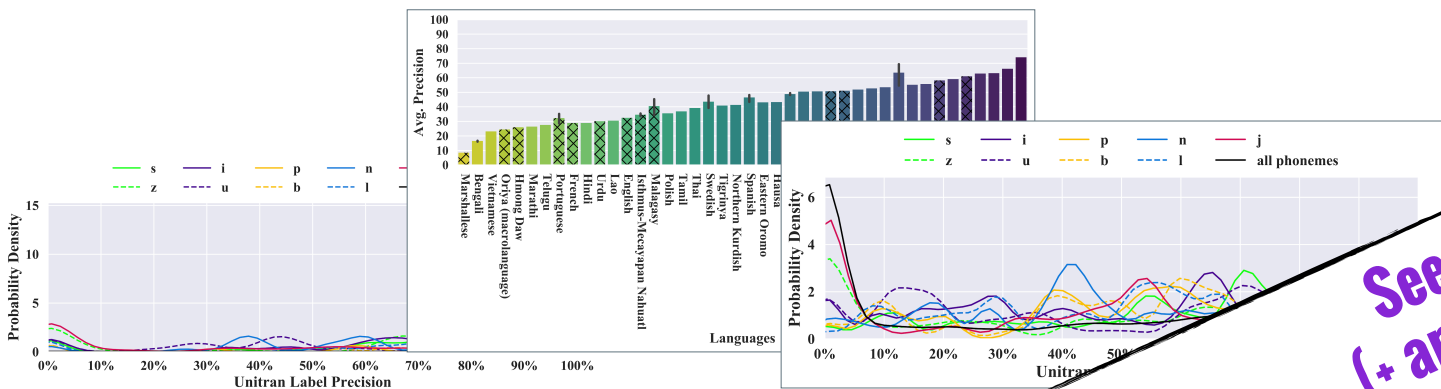**Spectral peak, COG, Duration, ...**

🤔 **Why provide both Unitran and High-Resource alignments?**

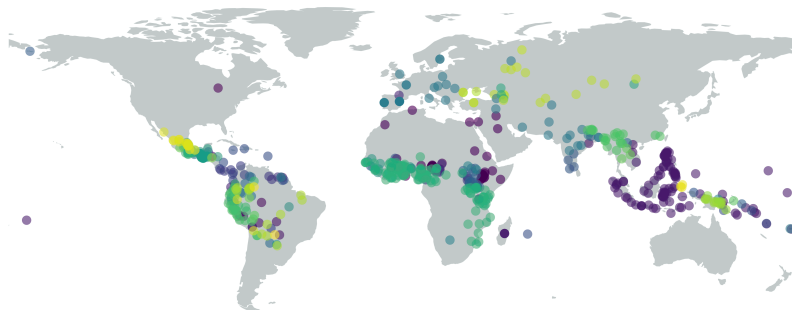Use multiple sets of alignments to assess **Unitran** alignment quality

- How much does quality vary across languages?

- Are certain phonemes more accurate than others?

- What about time alignment accuracy?



See paper!
(+ appendices)

# VoxClamantis v1.0 provides **tokens** of **phoneme-level measurements** in hundreds of languages!

- ‣ 690 recorded readings of the Bible
- ‣ 635 languages (ISO 639-3)
- ‣ 70 language families

- ‣ >400 million aligned phoneme-level segments
- ‣ Subsequent phonetic measures for all vowels and sibilants

# Case Studies

# Case studies with VoxClamantis v1.0

**Vowels**            **Sibilants**

~50 phonemes            /s/  /z/

**48 High-Resource Readings**

① **Reproduction of previous results validates resource**

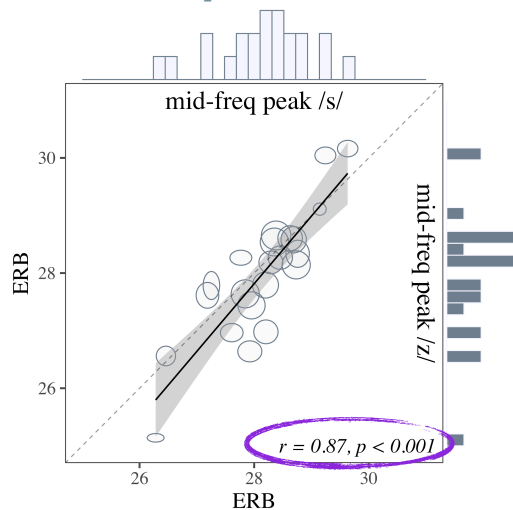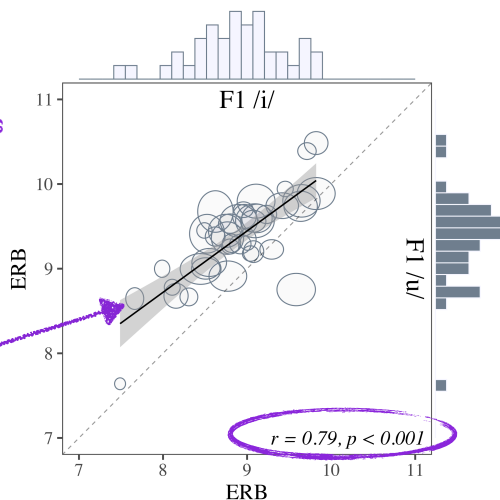② **Research at scale suggests general cross-linguistic principles**

# Are shared characteristics realized uniformly within languages?

(eg: vowel height, POA)     (eg: measures strongly correlated)

### *Formants*: **Vowels**

### *Mid-Freq Peak*: **Sibilants**



/i/, /u/: high vowels

F1 /i/
F1 /u/
ERB
ERB

*r = 0.79, p < 0.001*

(eg: language)

/s/, /z/: alveolar place of articulation

mid-freq peak /s/
mid-freq peak /z/
ERB
ERB

*r = 0.87, p < 0.001*

While variation exists across languages, within language F1 strongly correlated

## Reproduce previous results, but with many more languages
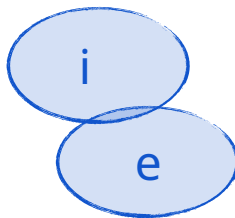
Supports hypothesis that this may be a universal principle

# Is inventory size correlated with articulatory precision?

**VOWELS**

**4 vowels** ⟶ **20 vowels**

i i:
ɪ
u u:
e
ɘ
ə
o
ɛ
ɜ:
ɔ: ɒ:
ɑ
æ
ɒ ɒ:
a:

**Marshallese** 🇲🇭                    **English** 🇬🇧

# Is inventory size correlated with articulatory precision?

**4 vowels** → **20 vowels**

Marshallese 🇲🇭

i
e
ɛ
æ

English 🇬🇧

i iː
ɪ
e
u uː
ɚ
ə ɝ
ɛ
ɜː
ɒ
æ
ɔ ɔː
aː
ɑ ɑː

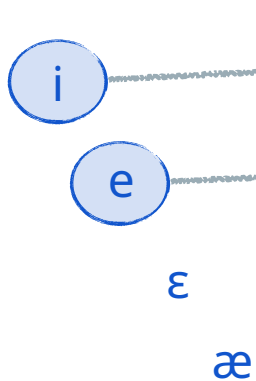**Is inventory size correlated with articulatory precision?**

**No**
(Spearman ρ = 0.11, p = 0.44;
Pearson r = 0.11, p = 0.46)

4 vowels ⟶ 20 vowels



Marshallese 🇲🇭     English 🇬🇧

**Previously shown,
but not possible to study at scale**

**Supports hypothesis
that this may [not] be a
universal principle**

27

CAUTION

B+     Utterance alignment     Filter -- in future, realign!

A- D+     Automatic phoneme labels     Better G(+A)2P

A 0% 😱     Alignment assessment!     Curate more resources!

B     Corpus representation
(e.g. speakers)     Curate more resources!

# Summary

voxclamantisproject.github.io

**VoxClamantis** v1.0 **corpus:**

/vɔks/ aligned phoneme-level segments in hundreds of languages
*57 high-resource, 690 first-pass*

😱 methodology is not perfect – version 1.0!
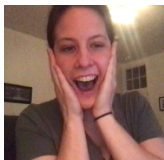
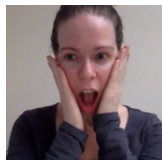⬇️ download 🥳 use for research ⬆️ contribute to v2.0!
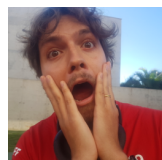
Questions!
Comments!
Contributions!

voxclamantisproject.github.io

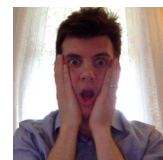voxclamantisproject@gmail.com

Elizabeth Salesky     Eleanor Chodroff     Tiago Pimentel     Matthew Wiesner
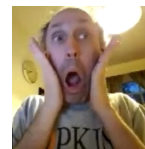
Ryan Cotterell     Alan W Black     Jason Eisner

VoxClamantis in deserto:
"a voice crying out in
the wilderness"