# Multilingual MFA: Forced Alignment on Low-Resource Related Languages
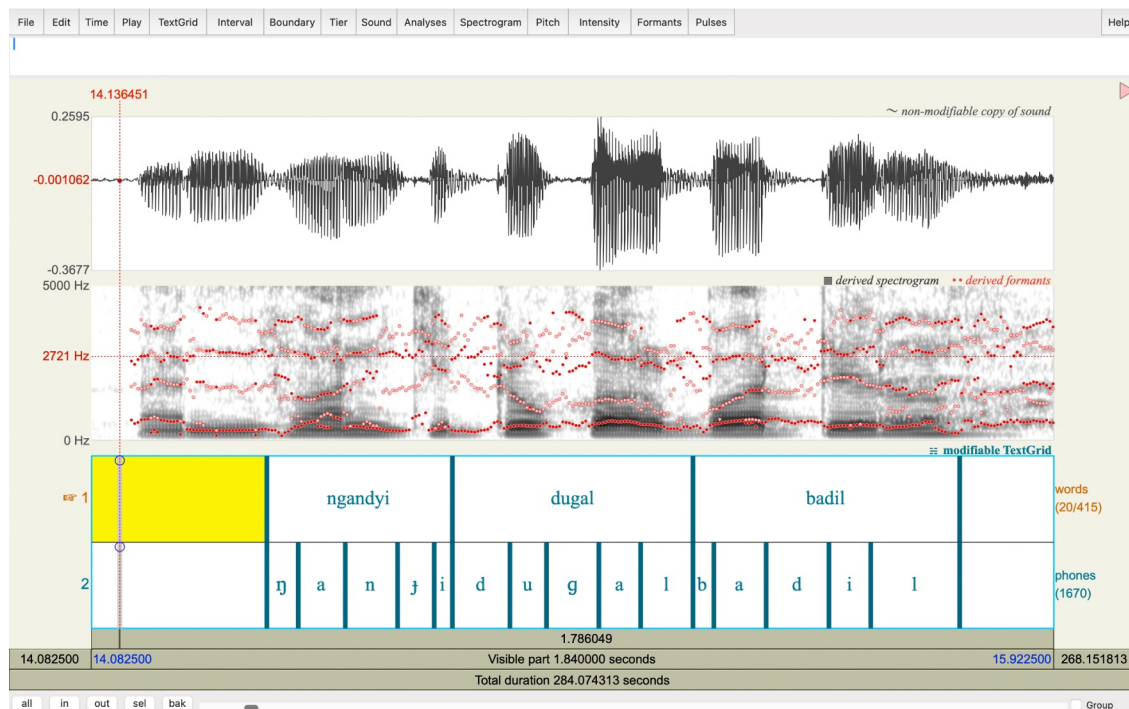
## Alessio Tosolini and Claire Bowern

Yale University

ICLDC / ComputEl 2025

1

# Roadmap

- Introduction : underresourced languages and forced alignment
- Models and methodology
- Results
    - Accuracy of the models
    - Analyses of the models
- Main takeaways

# Forced Alignment

Forced Alignment associates transcripts with audio and video (at utterance, word, or segment level).

It's incredibly useful for both linguistic research and community documentation projects.

Forced alignment requires an acoustic model and information about the grapheme to phoneme mappings (e.g. a dictionary of words and their phonemes). Acoustic model training is data hungry, and performance on languages across the world is very unequal.

McAuliffe et al 2017; Chodroff et al. 2024; DiCanio et al 2013; Babinski et al 2019

Various methods exist for increasing performance

Use high resource model (e.g. English)
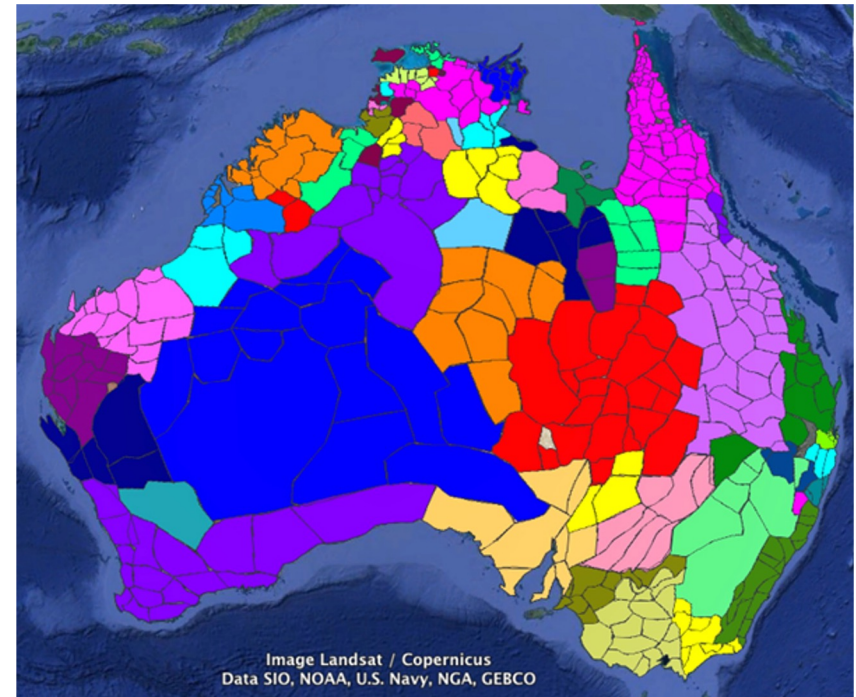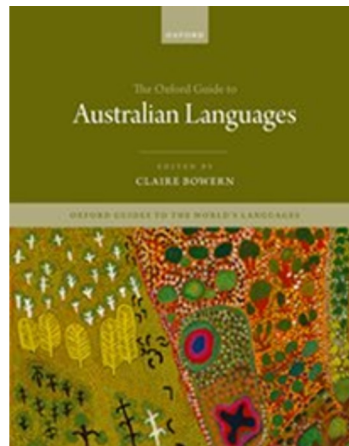
Adapt high resource model

Multilingual models

Here, we ask whether we can usefully combine language data from different languages from the same families (with very similar phoneme inventories) to get model improvement.

cf. San et al. 2021; Chodroff et al. 2024

4

# Australian Indigenous languages

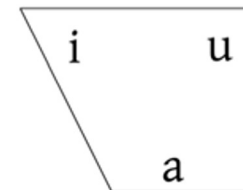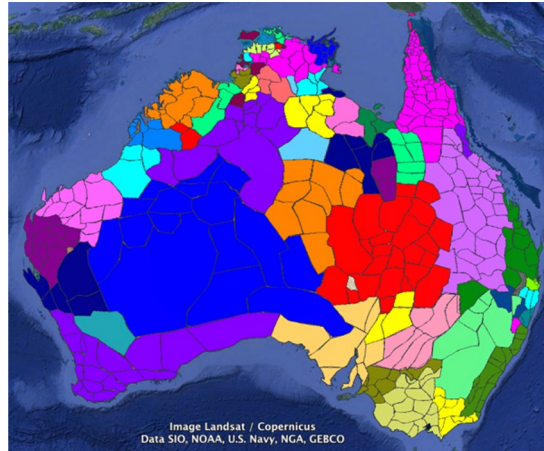**c. 440 languages**

**Similar phoneme inventories across the country**

<span style="color:darkred">**Chronically underrepresented in multilingual datasets (e.g. Common Voice)**</span>



The Oxford Guide to Australian Languages
EDITED BY CLAIRE BOWEN
OXFORD GUIDES TO THE WORLD'S LANGUAGES



Image Landsat / Copernicus
Data SIO, NOAA, U.S. Navy, NGA, GEBCO

|        | Labial | Lamino-dental | Apico-alveolar | Retroflex | Palatal | Velar |
|--------|--------|---------------|----------------|-----------|---------|-------|
| Nasal  | m      | n̪            | n              | ɳ         | ɲ       | ŋ     |
| Stop   | p      | t̪            | t              | ʈ         | c       | k     |
| Liquid |        | l̪            | l ɻ r          | ɭ         | ʎ       |       |
| Glide  | w      |               | j              |           |         |       |



i    u

a

# Can we take advantage of multilingual datasets to train better forced aligners?

# Models and methodology

Models:

- English adapted (base model trained on over 3000 hours of **global** data)
    - Yidiny
    - Big5 (Bardi, Gija, Ngaanyatjarra, Yan-nhangu, Yidiny)
    - Base model (McAuliffe and Sonderegger, 2024)
- From scratch
    - Yidiny
    - Big5 (Bardi, Gija, Ngaanyatjarra, Yan-nhangu, Yidiny)

Testing Datasets:

- Yidiny, seen data
- Yidiny, unseen data
- Kunbarlang

# Models and methodology

Models:

- English adapted (base model trained on over 3000 hours of **global** data)
    - Yidiny
    - Big5 (Bardi, Gija, Ngaanyatjarra, Yan-nhangu, Yidiny)
    - Base model (McAuliffe and Sonderegger, 2024)
- From scratch
    - Yidiny
    - Big5 (Bardi, Gija, Ngaanyatjarra, Yan-nhangu, Yidiny)

Testing Datasets:

- Yidiny, seen data
- Yidiny, unseen data
- Kunbarlang

| Language | Language Family | Reference | Collector | Minutes |
|---|---|---|---|---|
| Bardi | Nyulnyulan | A: Bowern_C05 | Claire Bowern | 108 |
| Gija | Jarrakan | E: 0098MDP0190 | Frances Kofod | 157 |
| Kunbarlang | Gunwinyguan | E: 0384SG0324 | Isabel O'Keefe; Ruth Singer | 16 |
| Ngaanyatjarra | Pama-Nyungan | P: WDVA1 | Inge Kral | 53 |
| Yan-nhangu | Pama-Nyungan | E: dk0046 | Claire Bowern | 290 |
| Yidiny | Pama-Nyungan | A: A2616 | R.M.W. Dixon | 50 |

Table 1: Corpus information. A: AIATSIS; E: Elar; P: Paradisec

# Evaluating an MFA model

- Method 1: Looking at how closely manual alignments match with MFA alignments
    - Can look at overall "accuracy"
    - Break it down by manner of articulation and place of articulation
        - Will English adapted models perform worse on e.g. nasals?

# Evaluating an MFA model

- Method 1: Looking at how closely manual alignments match with MFA alignments
  - Can look at overall "accuracy"
  - Break it down by manner of articulation and place of articulation
    - Will English adapted models perform worse on e.g. nasals?
- Method 2: Looking at how closely manual *analyses* match with MFA alignments
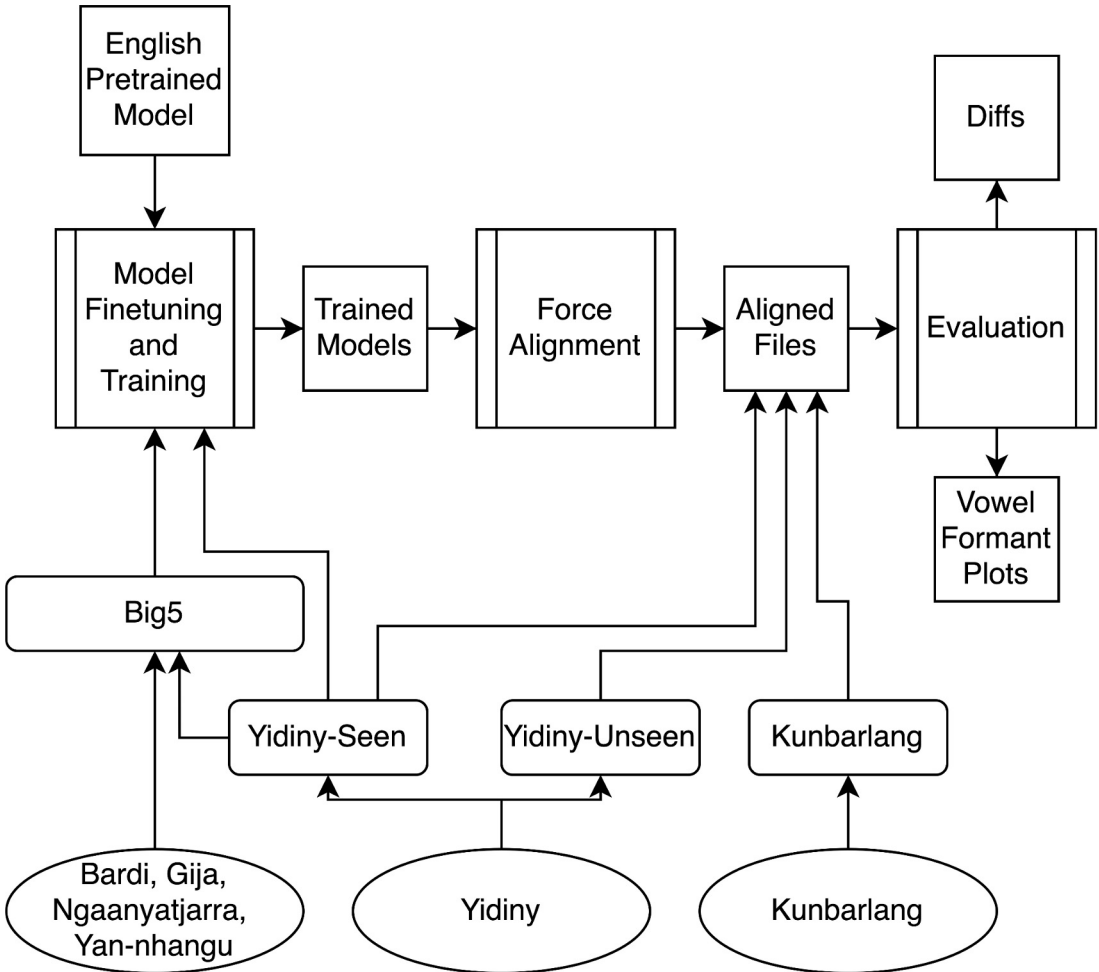  - We look at vowel formant plots

# Evaluating an MFA model

- Method 1: Looking at how closely manual alignments match with MFA alignments
  - Can look at overall "accuracy"
  - Break it down by manner of articulation and place of articulation
    - Will English adapted models perform worse on e.g. nasals?
- Method 2: Looking at how closely manual *analyses* match with MFA alignments
  - We look at vowel formant plots
- Human analyses vary! There is no one "gold standard"

# Pipeline

# Pipeline



Models we are comparing

Datasets

English Pretrained Model

Model Finetuning and Training → Trained Models → Force Alignment → Aligned Files → Evaluation

Diffs

Vowel Formant Plots

Big5

Yidiny-Seen

Yidiny-Unseen

Kunbarlang

Bardi, Gija, Ngaanyatjarra, Yan-nhangu

Yidiny

Kunbarlang

13

# Results: Accuracy

Things to note:
- They ideally should all have approx. **mean** 0 - which we see!
- The more spread (**sd**) a plot has, the more variation the forced aligner produces





Figure 1: Onset boundary differences for all models across all testing settings.

# Results: Accuracy

Things to note:
- They ideally should all have approx. **mean** 0 - which we see!
- The more spread (**sd**) a plot has, the more variation the forced aligner produces



Figure 1: Onset boundary differences for all models across all testing settings.

# Results: Accuracy

Things to note:
- They ideally should all have approx. **mean** 0 - which we see!
- The more spread (**sd**) a plot has, the more variation the forced aligner produces



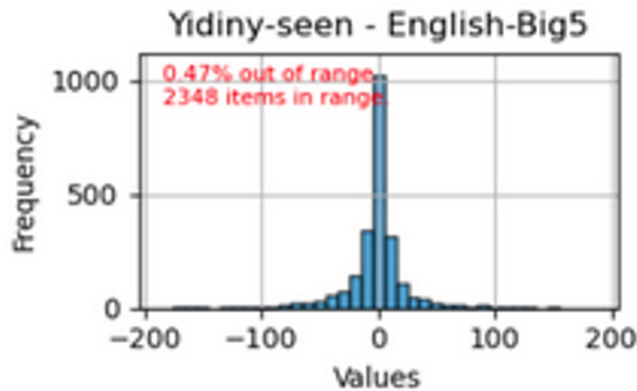Figure 1: Onset boundary differences for all models across all testing settings.

# Results: Accuracy

Things to note:
- They ideally should all have approx. **mean** 0 - which we see!
- The more spread (**sd**) a plot has, the more variation the forced aligner produces
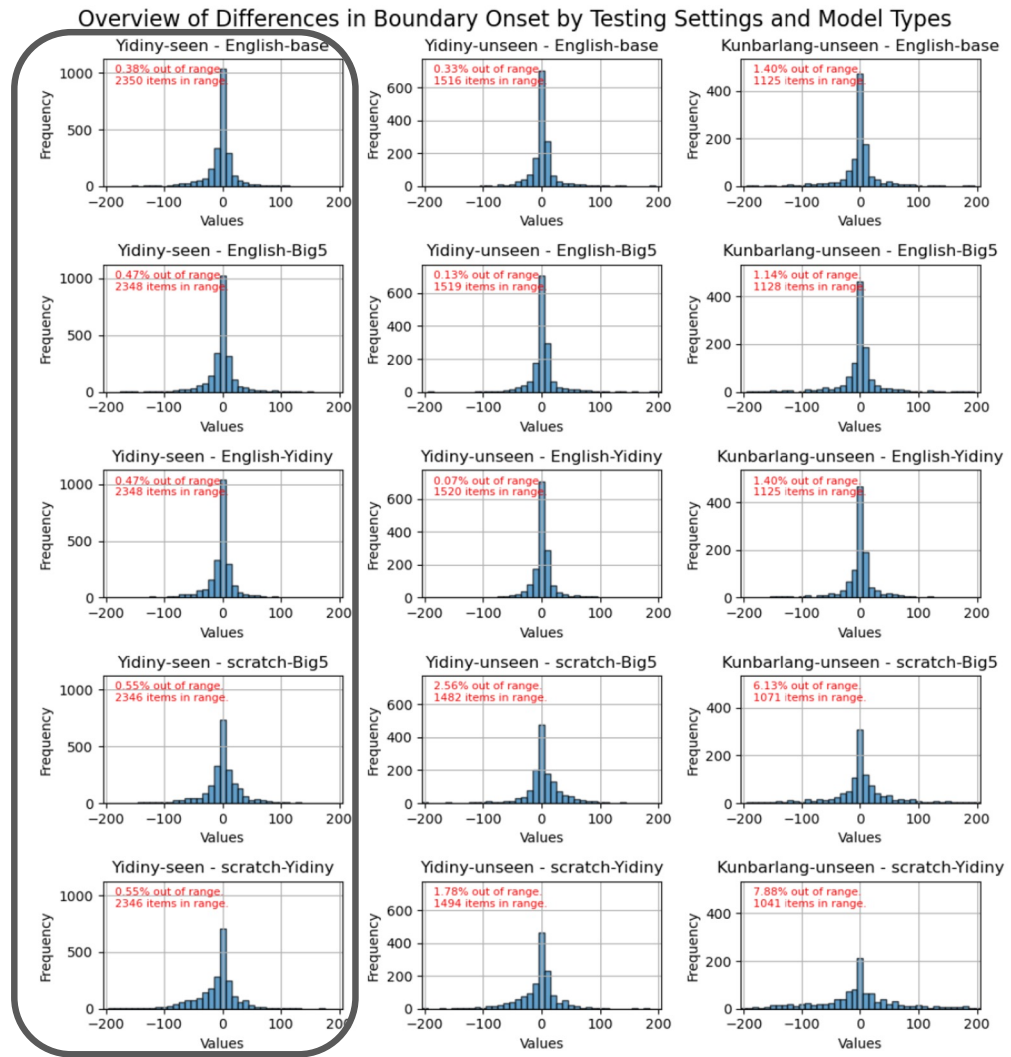


Figure 1: Onset boundary differences for all models across all testing settings.

# Results: Accuracy - Yidiny seen testing setting



Heatmap of onset boundary mean differences - Yidiny-seen

# Results: Accuracy - Yidiny seen testing setting



Heatmap of onset boundary mean differences - Yidiny-seen

|  | All Vowels | Short Vowels | Long Vowels | All Stops | Nasals | Rhotics | Trill | Approximants | Overall diff. | Overall abs. diff. |
|---|---|---|---|---|---|---|---|---|---|---|
| English-Big5 | -1.21 | -0.90 | -3.19 | -1.83 | -1.68 | -7.42 | -13.43 | -7.23 | -2.23 | 16.38 |
| English-Yidiny | -2.06 | -1.56 | -5.19 | -2.42 | -1.48 | -6.82 | -19.99 | -6.92 | -2.80 | 16.70 |
| English-base | -3.31 | -2.80 | -6.53 | -2.62 | -2.23 | -7.73 | -18.12 | -7.60 | -3.73 | 17.04 |
| scratch-Big5 | 4.32 | 2.32 | 16.95 | -2.54 | 5.21 | -12.27 | -26.87 | -16.87 | 0.59 | 23.92 |
| scratch-Yidiny | -0.22 | -0.51 | 1.62 | -10.90 | -16.10 | -30.60 | -39.06 | -31.55 | -8.21 | 24.94 |

# Results: Accuracy - Yidiny seen testing setting



Heatmap of onset boundary mean differences - Yidiny-seen

# Results: Accuracy - Yidiny seen testing setting



Heatmap of onset boundary mean differences - Yidiny-seen

# Results: Accuracy - Yidiny seen testing setting



Heatmap of onset boundary mean differences - Yidiny-seen

| | All Vowels | Short Vowels | Long Vowels | All Stops | Nasals | Rhotics | Trill | Approximants | Overall diff. | Overall abs. diff. |
|---|---|---|---|---|---|---|---|---|---|---|
| English-Big5 | -1.21 | -0.90 | -3.19 | -1.83 | -1.68 | -7.42 | -13.43 | -7.23 | -2.23 | 16.38 |
| English-Yidiny | -2.06 | -1.56 | -5.19 | -2.42 | -1.48 | -6.82 | -19.99 | -6.92 | -2.80 | 16.70 |
| English-base | -3.31 | -2.80 | -6.53 | -2.62 | -2.23 | -7.73 | -18.12 | -7.60 | -3.73 | 17.04 |
| scratch-Big5 | 4.32 | 2.32 | 16.95 | -2.54 | 5.21 | -12.27 | -26.87 | -16.87 | 0.59 | 23.92 |
| scratch-Yidiny | -0.22 | -0.51 | 1.62 | -10.90 | -16.10 | -30.60 | -39.06 | -31.55 | -8.21 | 24.94 |

# Results: Precision - Yidiny seen testing setting



Heatmap of onset boundary standard deviations - Yidiny-seen

| | All Vowels | Short Vowels | Long Vowels | All Stops | Nasals | Rhotics | Trill | Approximants | Overall diff. |
|---|---|---|---|---|---|---|---|---|---|
| English-Big5 | 36.96 | 33.85 | 52.45 | 44.83 | 32.61 | 34.65 | 33.30 | 34.17 | 37.43 |
| English-Yidiny | 36.94 | 33.38 | 54.18 | 44.50 | 33.69 | 42.30 | 32.70 | 37.38 | 37.96 |
| English-base | 38.62 | 35.55 | 54.06 | 44.95 | 32.17 | 39.03 | 30.66 | 34.79 | 38.47 |
| scratch-Big5 | 45.04 | 40.80 | 64.36 | 46.63 | 40.10 | 36.35 | 42.79 | 46.38 | 45.22 |
| scratch-Yidiny | 37.61 | 38.10 | 34.30 | 40.30 | 43.38 | 38.12 | 40.88 | 49.62 | 41.76 |

# Results: Precision - Yidiny seen testing setting



Heatmap of onset boundary standard deviations - Yidiny-seen

| | All Vowels | Short Vowels | Long Vowels | All Stops | Nasals | Rhotics | Trill | Approximants | Overall diff. |
|---|---|---|---|---|---|---|---|---|---|
| English-Big5 | 36.96 | 33.85 | 52.45 | 44.83 | 32.61 | 34.65 | 33.30 | 34.17 | 37.43 |
| English-Yidiny | 36.94 | 33.38 | 54.18 | 44.50 | 33.69 | 42.30 | 32.70 | 37.38 | 37.96 |
| English-base | 38.62 | 35.55 | 54.06 | 44.95 | 32.17 | 39.03 | 30.66 | 34.79 | 38.47 |
| scratch-Big5 | 45.04 | 40.80 | 64.36 | 46.63 | 40.10 | 36.35 | 42.79 | 46.38 | 45.22 |
| scratch-Yidiny | 37.61 | 38.10 | 34.30 | 40.30 | 43.38 | 38.12 | 40.88 | 49.62 | 41.76 |

# Results: Precision - Yidiny seen testing setting



Heatmap of onset boundary standard deviations - Yidiny-seen

# Precision Yidiny-seen Summary

In the seen language seen data setting:
- Multilingual models show (slightly) higher precision than monolingual models
- Models trained from scratch are <u>slightly</u> less precise and less accurate than English-based models



Heatmap of onset boundary mean differences - Yidiny-seen



Heatmap of onset boundary standard deviations - Yidiny-seen

# Results: Accuracy - Yidiny unseen testing setting



Heatmap of onset boundary mean differences - Yidiny-unseen

Heatmap of onset boundary standard deviations - Yidiny-unseen

# Results: Accuracy - Yidiny unseen testing setting



Heatmap of onset boundary mean differences - Yidiny-unseen

Heatmap of onset boundary standard deviations - Yidiny-unseen

# Results: Accuracy - Yidiny unseen testing setting



Heatmap of onset boundary mean differences - Yidiny-unseen

Heatmap of onset boundary standard deviations - Yidiny-unseen

# Results: Accuracy - Yidiny unseen testing setting



Heatmap of onset boundary mean differences - Yidiny-unseen
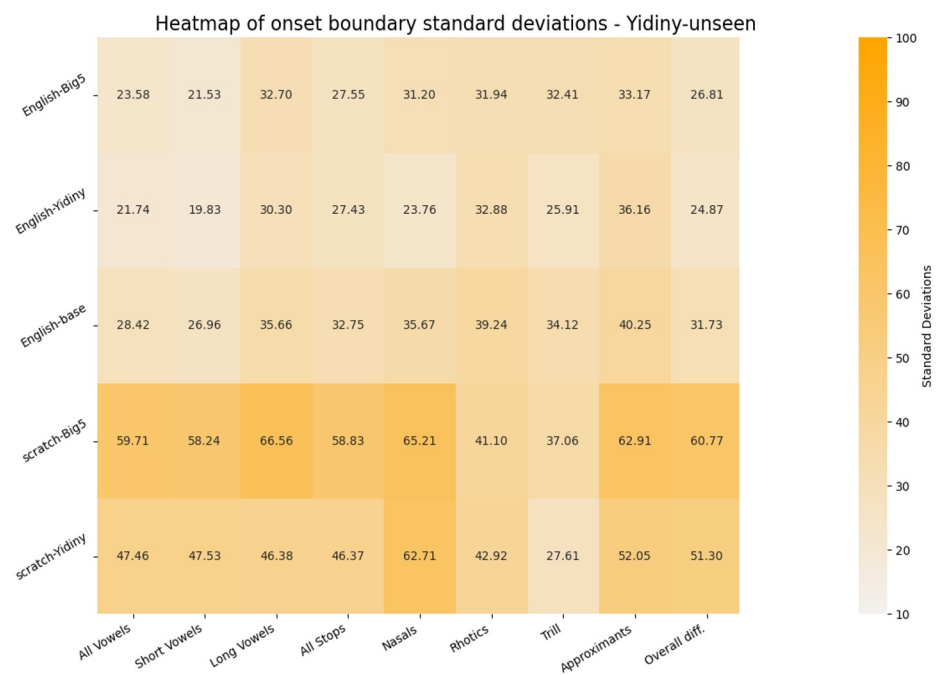
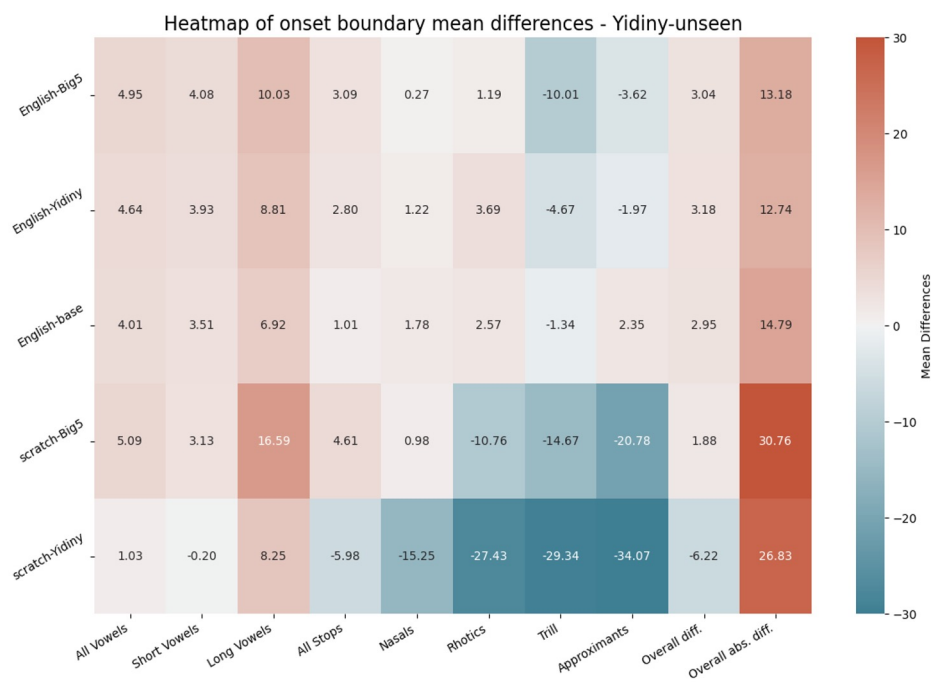Heatmap of onset boundary standard deviations - Yidiny-unseen

# Precision Yidiny-unseen Summary

In the seen language unseen data setting:
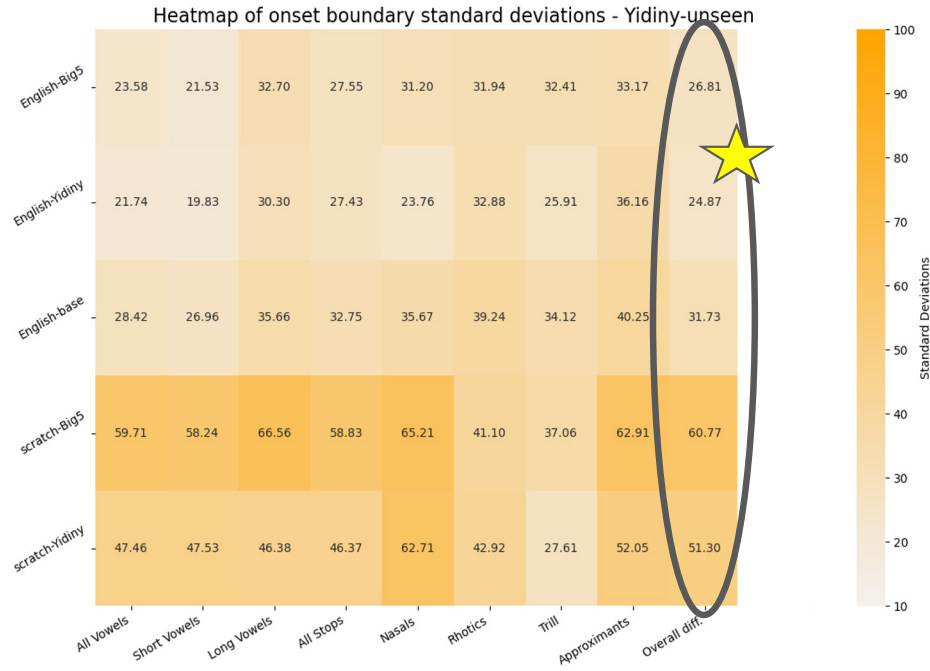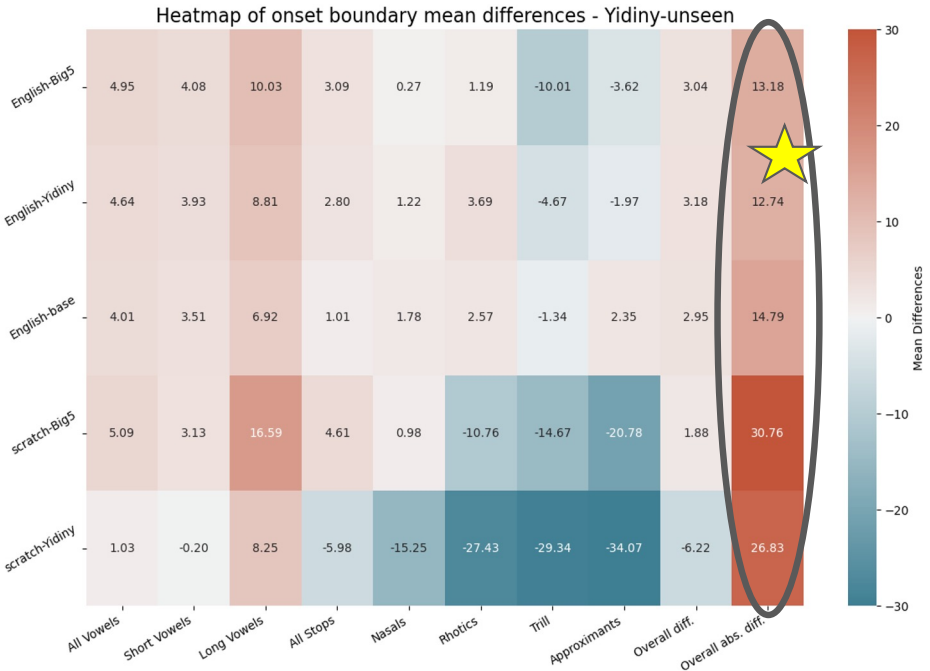- Multilingual models show (slightly) lower precision than monolingual models
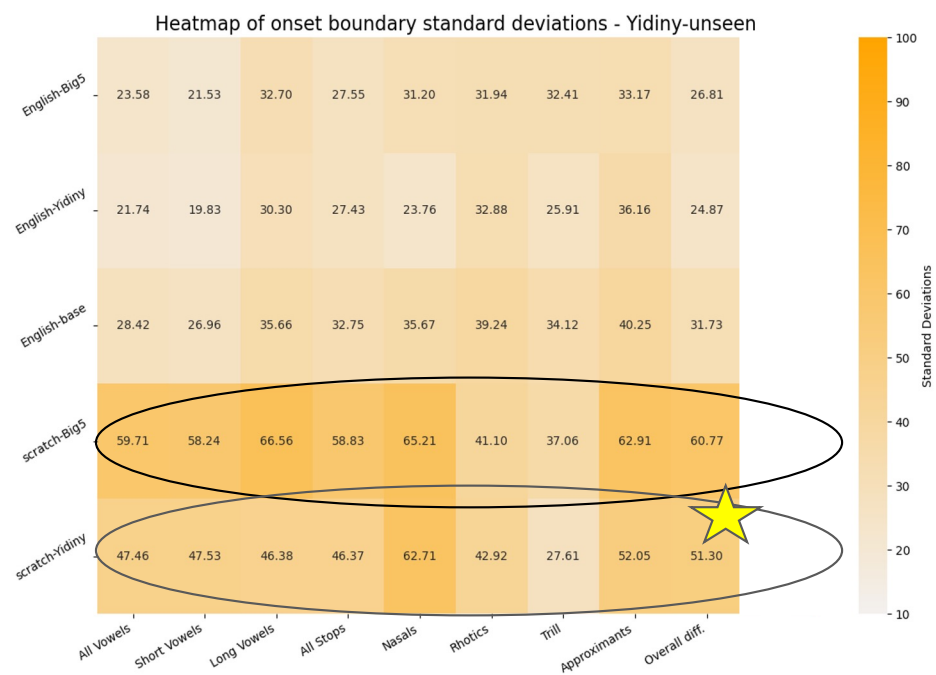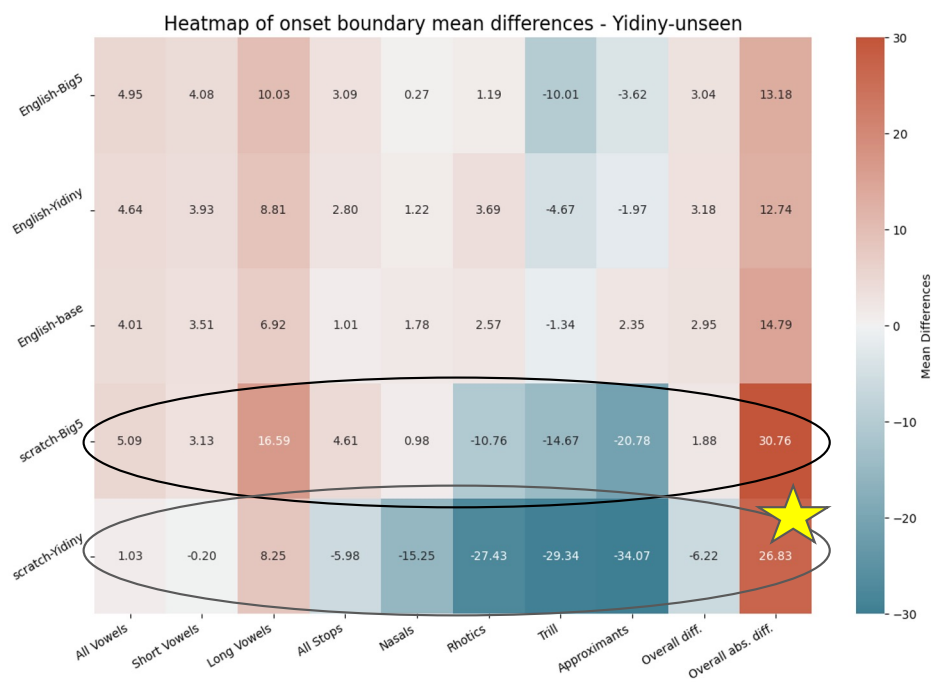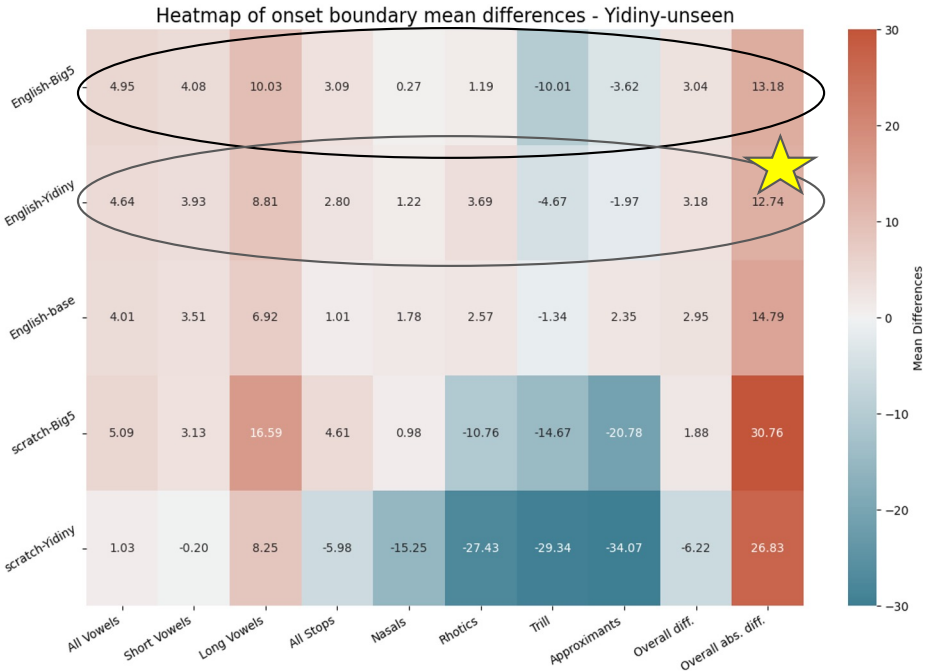- Models trained from scratch are quite less precise and less accurate than English-based models

### Heatmap of onset boundary mean differences - Yidiny-unseen

| | All Vowels | Short Vowels | Long Vowels | All Stops | Nasals | Rhotics | Trill | Approximants | Overall diff. | Overall abs. diff. |
|---|---|---|---|---|---|---|---|---|---|---|
| English-Big5 | 4.95 | 4.08 | 10.03 | 3.09 | 0.27 | 1.19 | -10.01 | -3.62 | 3.04 | 13.18 |
| English-Yidiny | 4.64 | 3.93 | 8.81 | 2.80 | 1.22 | 3.69 | -4.67 | -1.97 | 3.18 | 12.74 |
| English-base | 4.01 | 3.51 | 6.92 | 1.01 | 1.78 | 2.57 | -1.34 | 2.35 | 2.95 | 14.79 |
| scratch-Big5 | 5.09 | 3.13 | 16.59 | 4.61 | 0.98 | -10.76 | -14.67 | -20.78 | 1.88 | 30.76 |
| scratch-Yidiny | 1.03 | -0.20 | 8.25 | -5.98 | -15.25 | -27.43 | -29.34 | -34.07 | -6.22 | 26.83 |

### Heatmap of onset boundary standard deviations - Yidiny-unseen

| | All Vowels | Short Vowels | Long Vowels | All Stops | Nasals | Rhotics | Trill | Approximants | Overall diff. |
|---|---|---|---|---|---|---|---|---|---|
| English-Big5 | 23.58 | 21.53 | 32.70 | 27.55 | 31.20 | 31.94 | 32.41 | 33.17 | 26.81 |
| English-Yidiny | 21.74 | 19.83 | 30.30 | 27.43 | 23.76 | 32.88 | 25.91 | 36.16 | 24.87 |
| English-base | 28.42 | 26.96 | 35.66 | 32.75 | 35.67 | 39.24 | 34.12 | 40.25 | 31.73 |
| scratch-Big5 | 59.71 | 58.24 | 66.56 | 58.83 | 65.21 | 41.10 | 37.06 | 62.91 | 60.77 |
| scratch-Yidiny | 47.46 | 47.53 | 46.38 | 46.37 | 62.71 | 42.92 | 27.61 | 52.05 | 51.30 |

31

# Results: Accuracy - Kunbarlang (unseen) testing setting



Heatmap of onset boundary mean differences - Kunbarlang-unseen

Heatmap of onset boundary standard deviations - Kunbarlang-unseen

# Results: Accuracy - Kunbarlang (unseen) testing setting



Heatmap of onset boundary mean differences - Kunbarlang-unseen

| | All Vowels | All Stops | Nasals | Rhotics | Approximants | Overall diff. | Overall abs. diff. |
|---|---|---|---|---|---|---|---|
| English-Big5 | -2.08 | -5.06 | -0.59 | 14.33 | 5.21 | -0.63 | 23.44 |
| English-Yidiny | -0.05 | -4.83 | -0.07 | 12.44 | 9.63 | 0.74 | 24.40 |
| English-base | 0.42 | -3.58 | 1.30 | 19.83 | 13.61 | 2.38 | 25.20 |
| scratch-Big5 | 22.65 | 5.89 | 14.53 | -16.54 | -3.85 | 11.88 | 50.43 |
| scratch-Yidiny | 13.01 | -3.67 | 3.13 | -21.35 | -50.90 | -2.18 | 70.95 |

Heatmap of onset boundary standard deviations - Kunbarlang-unseen

| | All Vowels | All Stops | Nasals | Rhotics | Approximants | Overall diff. |
|---|---|---|---|---|---|---|
| English-Big5 | 54.40 | 40.07 | 53.00 | 83.46 | 45.53 | 52.73 |
| English-Yidiny | 60.14 | 42.97 | 54.17 | 78.17 | 66.30 | 57.81 |
| English-base | 60.30 | 39.64 | 53.63 | 85.43 | 65.79 | 57.92 |
| scratch-Big5 | 96.59 | 117.58 | 86.78 | 71.61 | 86.05 | 98.04 |
| scratch-Yidiny | 113.29 | 134.80 | 104.79 | 126.70 | 92.23 | 117.06 |

# Results: Accuracy - Kunbarlang (unseen) testing setting



Heatmap of onset boundary mean differences - Kunbarlang-unseen

Heatmap of onset boundary standard deviations - Kunbarlang-unseen

# Results: Accuracy - Kunbarlang (unseen) testing setting



Heatmap of onset boundary mean differences - Kunbarlang-unseen

|  | All Vowels | All Stops | Nasals | Rhotics | Approximants | Overall diff. | Overall abs. diff. |
|---|---|---|---|---|---|---|---|
| English-Big5 | -2.08 | -5.06 | -0.59 | 14.33 | 5.21 | -0.63 | 23.44 |
| English-Yidiny | -0.05 | -4.83 | -0.07 | 12.44 | 9.63 | 0.74 | 24.40 |
| English-base | 0.42 | -3.58 | 1.30 | 19.83 | 13.61 | 2.38 | 25.20 |
| scratch-Big5 | 22.65 | 5.89 | 14.53 | -16.54 | -3.85 | 11.88 | 50.43 |
| scratch-Yidiny | 13.01 | -3.67 | 3.13 | -21.35 | -50.90 | -2.18 | 70.95 |

Heatmap of onset boundary standard deviations - Kunbarlang-unseen

|  | All Vowels | All Stops | Nasals | Rhotics | Approximants | Overall diff. |
|---|---|---|---|---|---|---|
| English-Big5 | 54.40 | 40.07 | 53.00 | 83.46 | 45.53 | 52.73 |
| English-Yidiny | 60.14 | 42.97 | 54.17 | 78.17 | 66.30 | 57.81 |
| English-base | 60.30 | 39.64 | 53.63 | 85.43 | 65.79 | 57.92 |
| scratch-Big5 | 96.59 | 117.58 | 86.78 | 71.61 | 86.05 | 98.04 |
| scratch-Yidiny | 113.29 | 134.80 | 104.79 | 126.70 | 92.23 | 117.06 |

35

# Precision Kunbarlang Summary

In the unseen language seen data setting:
- Multilingual models show more precision than monolingual models
- Models trained from scratch are <u>much</u> less precise and less accurate than English-based models



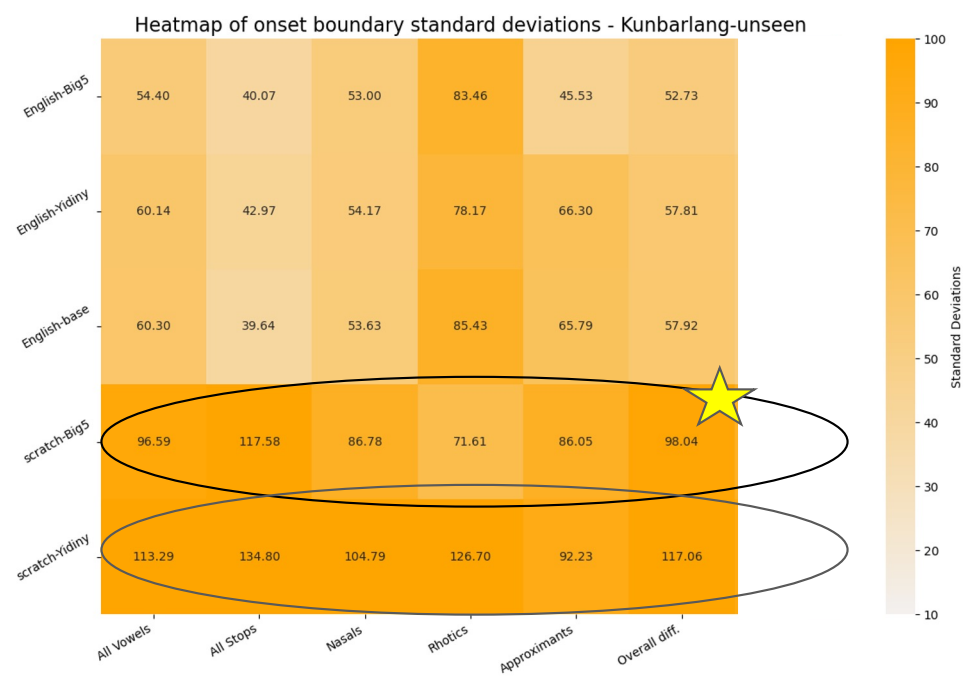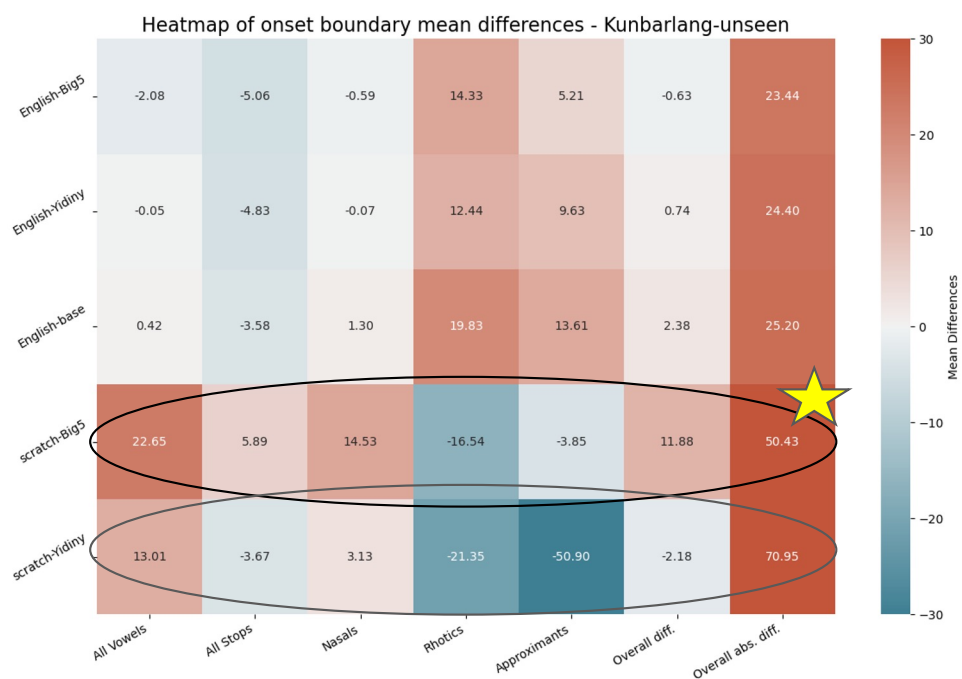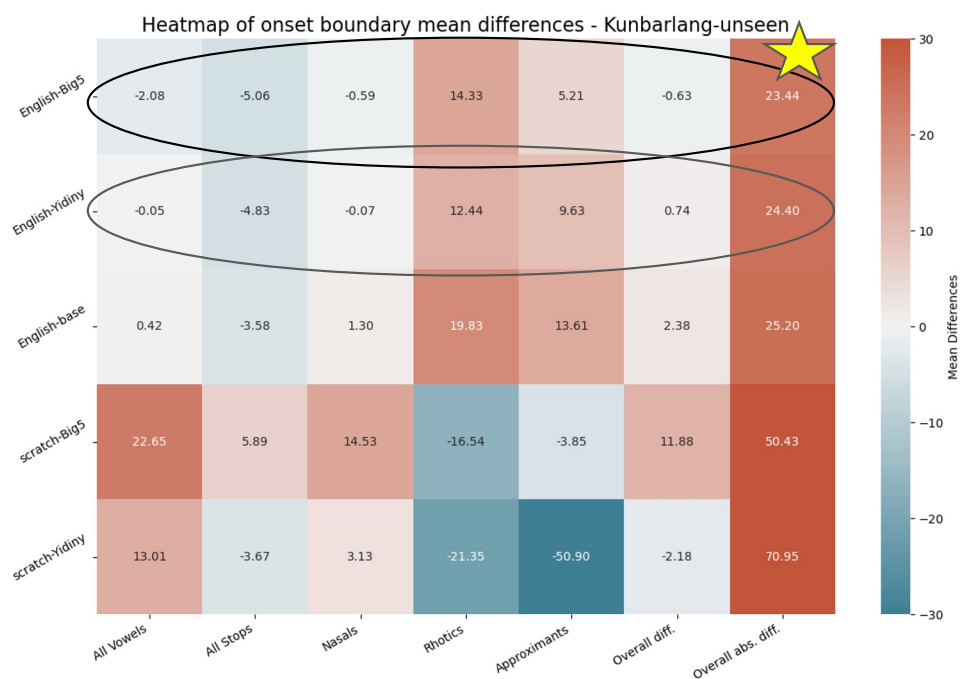Heatmap of onset boundary mean differences - Kunbarlang-unseen

|  | All Vowels | All Stops | Nasals | Rhotics | Approximants | Overall diff. | Overall abs. diff. |
|---|---|---|---|---|---|---|---|
| English-BigS | -2.08 | -5.06 | -0.59 | 14.33 | 5.21 | -0.63 | 23.44 |
| English-Yidiny | -0.05 | -4.83 | -0.07 | 12.44 | 9.63 | 0.74 | 24.40 |
| English-base | 0.42 | -3.58 | 1.30 | 19.83 | 13.61 | 2.38 | 25.20 |
| scratch-BigS | 22.65 | 5.89 | 14.53 | -16.54 | -3.85 | 11.88 | 50.43 |
| scratch-Yidiny | 13.01 | -3.67 | 3.13 | -21.35 | -50.90 | -2.18 | 70.95 |



Heatmap of onset boundary standard deviations - Kunbarlang-unseen

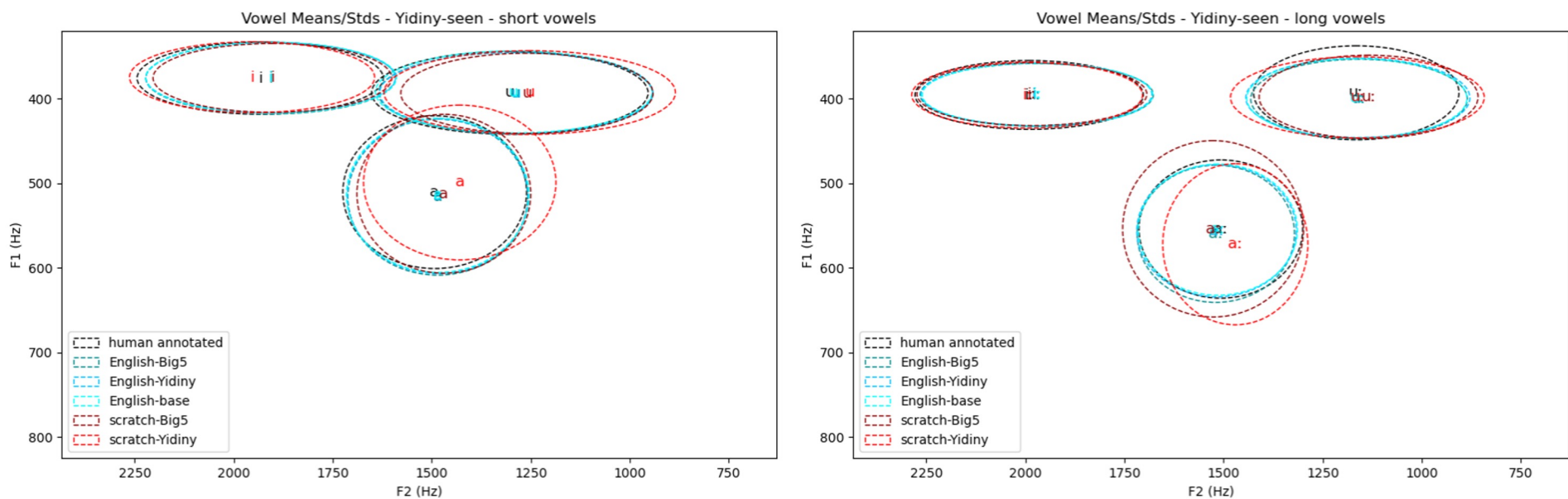|  | All Vowels | All Stops | Nasals | Rhotics | Approximants | Overall diff. |
|---|---|---|---|---|---|---|
| English-BigS | 54.40 | 40.07 | 53.00 | 83.46 | 45.53 | 52.73 |
| English-Yidiny | 60.14 | 42.97 | 54.17 | 78.17 | 66.30 | 57.81 |
| English-base | 60.30 | 39.64 | 53.63 | 85.43 | 65.79 | 57.92 |
| scratch-BigS | 96.59 | 117.58 | 86.78 | 71.61 | 86.05 | 98.04 |
| scratch-Yidiny | 113.29 | 134.80 | 104.79 | 126.70 | 92.23 | 117.06 |

# Precision and Accuracy Summary

Across all settings:
- Precision and accuracy seem highly correlated
- Multilingual training data:
    - improved performance in the Yidiny-seen and Kunbarlang-unseen settings
    - slightly decreased performance in the Yidiny-unseen setting
- Models trained from scratch consistently perform worse than English-based models
    - The performance gap was Kunbarlang >> Yidiny-unseen >> Yidiny-seen
    - In the Kunbarlang setting, multilingual training data had the biggest (positive) impact

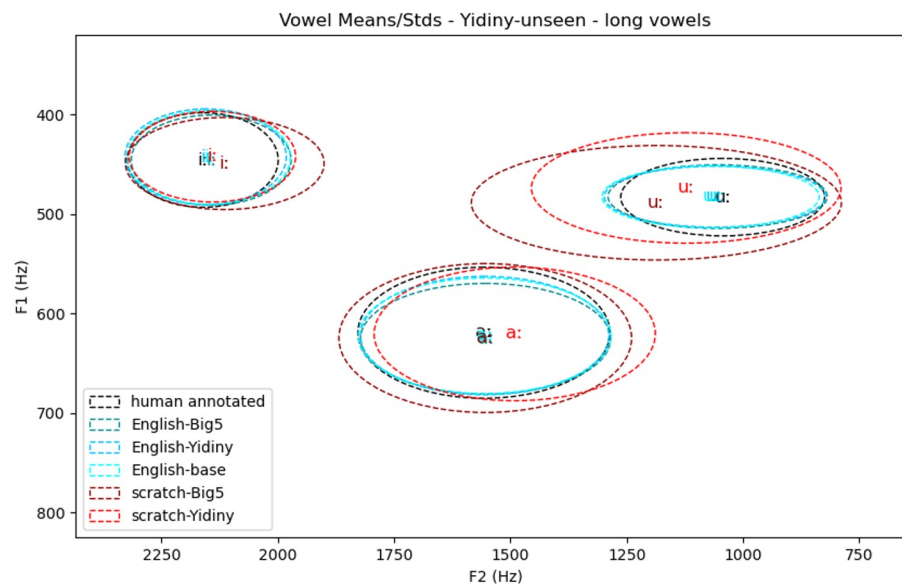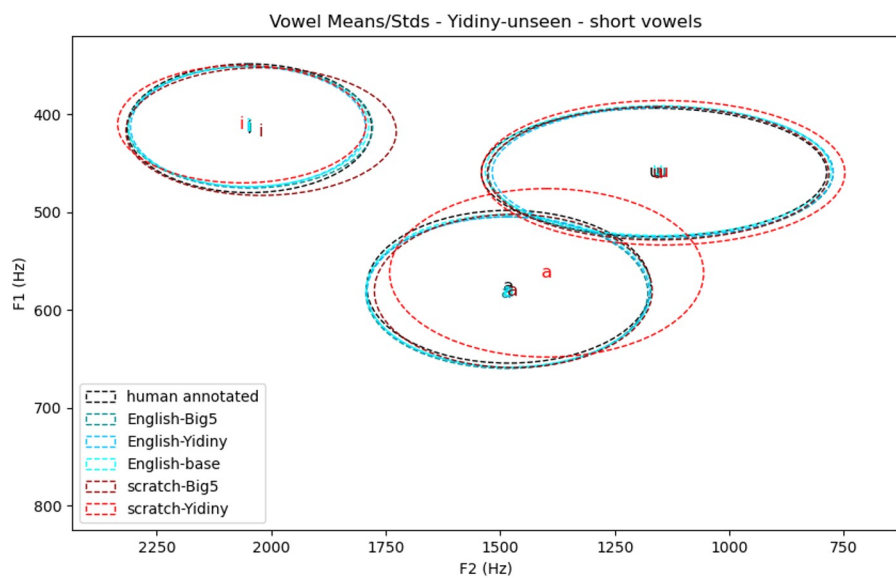# Comparison of Analyses: Yidiny seen



Vowel Means/Stds - Yidiny-seen - short vowels

Vowel Means/Stds - Yidiny-seen - long vowels

All analyses are similar, except the scratch-Yidiny model.
Long vowels show more variation for models trained from scratch.
**Size of ellipses is large because of difficulty in formant extraction (noisy data)**
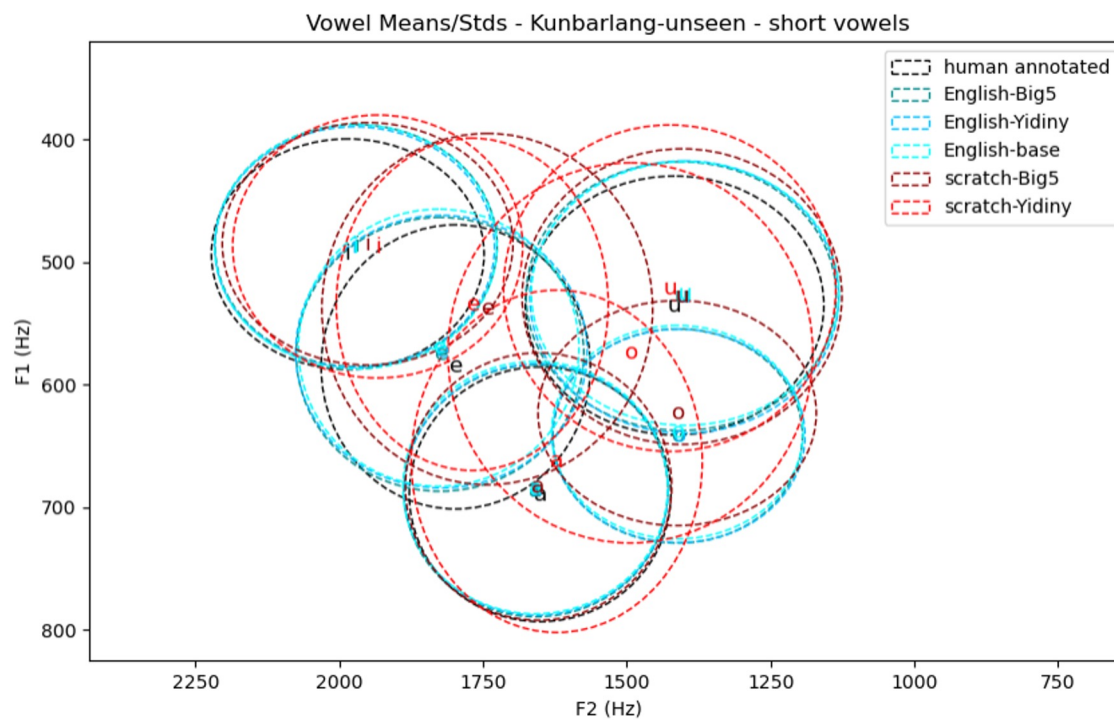
# Comparison of Analyses: Yidiny unseen



All analyses are similar, except the scratch-Yidiny model.

Long vowels show more variation for models trained from scratch.

# Comparison of Analyses: Kunbarlang (unseen)

Models from scratch:

- Struggle with /e o/
  - /e/ absent in all training data, /o/ only present, but rare in one language of Big5



Vowel Means/Stds - Kunbarlang-unseen - short vowels

Legend:
- human annotated
- English-Big5
- English-Yidiny
- English-base
- scratch-Big5
- scratch-Yidiny

# Comparison of Analyses Summary

In the Yidiny-seen and Yidiny-unseen settings:
- All models, except the model trained from scratch on only Yidiny data, gave approximately the same analyses
- Multilingual models provided better analyses

In the Kunbarlang setting:
- English-based models gave analyses nearly identical to the manually annotated analysis
- Models trained from scratch struggle with phones not in the training data (the mid-vowels)

In all settings:
- English-based models > scratch-Big5 >> scratch-Yidiny

# Main takeaways

- Which models work best?
    - English based models work best. Even the off-the-shelf English model was better than all the models trained from scratch in almost all ways
    - Models trained from scratch were *almost* as good as the English models in the Yidiny-seen setting
- Does multilingual training data improve MFA?
    - Yes!
        - Slight improvements for English models
            - Biggest improvements by natural class came from natural classes not in the training data (e.g. trills for English)
        - Larger improvements for models trained from scratch
- How do analyses compare?
    - English based models ≈ manual annotation
    - The multilingual model trained from scratch also ≈ manual annotation for tokens it has trained extensively on

# Low-Resource Forced Alignment and Future work

Low-Resource Forced Alignment:

- The Yidiny-seen setting is thus most similar to real settings
    - Smallest difference between the English-based and from-scratch models
- All models worked quite well in the Yidiny-seen setting
- **Recommendation: use an adapted English based model if you have less than 30 minutes of training data**

Future work:

- Optimizing forced alignment for settings like "Yidiny-seen"
- Hyperparameter tuning
- Data augmentation

# References + Contact

Chodroff, Eleanor, E. Ahn, and Hossep Dolatian. 2024. Comparing language-specific and cross-language acoustic models for low-resource phonetic forced alignment. Language Documentation & Conservation

Magueresse, A., Carles, V., & Heetderks, E. (2020). Low-resource Languages: A Review of Past Work and Future Challenges. *arXiv [Cs.CL]*. Retrieved from http://arxiv.org/abs/2006.07264

McAuliffe, Michael, and Morgan Sonderegger. 2024. English mfa acoustic model v3.1.0.Technical report, https://mfa-models.readthedocs.io/acoustic/English/EnglishMFAacousticmodelv3_1_0.html.

McAuliffe, Michael, Michaela Socolof, Sarah Mihuc, Michael Wagner, and Morgan Sonderegger. 2017.  Montreal forced aligner: Trainable text-speech alignment using kaldi. In Interspeech, volume 2017, pages 498–502.

Wu, H., Yun, J., Li, X. *et al.* Using a forced aligner for prosody research. *Humanit Soc Sci Commun* **10**, 429 (2023). https://doi.org/10.1057/s41599-023-01931-4

Young, N. J. and M. McGarrah, "Forced alignment for Nordic languages: Rapidly constructing a high-quality prototype," *Nordic Journal of Linguistics*, vol. 46, no. 1, pp. 105–131, May 2023.

alessio.tosolini@yale.edu                                   claire.bowern@yale.edu

# Acknowledgements

- ELAR, Paradisec, and AIATSIS archives
- Audiences at Yale and SYNC for feedback
- The elders and others who worked with linguists to make records of their languages
- National Science Foundation's Linguistics and DEL-DLI programs for long-term funding of research on language and with language communities, particularly in their support for broader impacts (earlier work funded by BCS-0844550, BCS-1423711, and BCS-2116164)

# Appendix A: PoA Accuracy Yidiny seen



Heatmap of onset boundary mean differences - Yidiny-seen

Heatmap of onset boundary standard deviations - Yidiny-seen

# Appendix A: PoA Accuracy Yidiny unseen



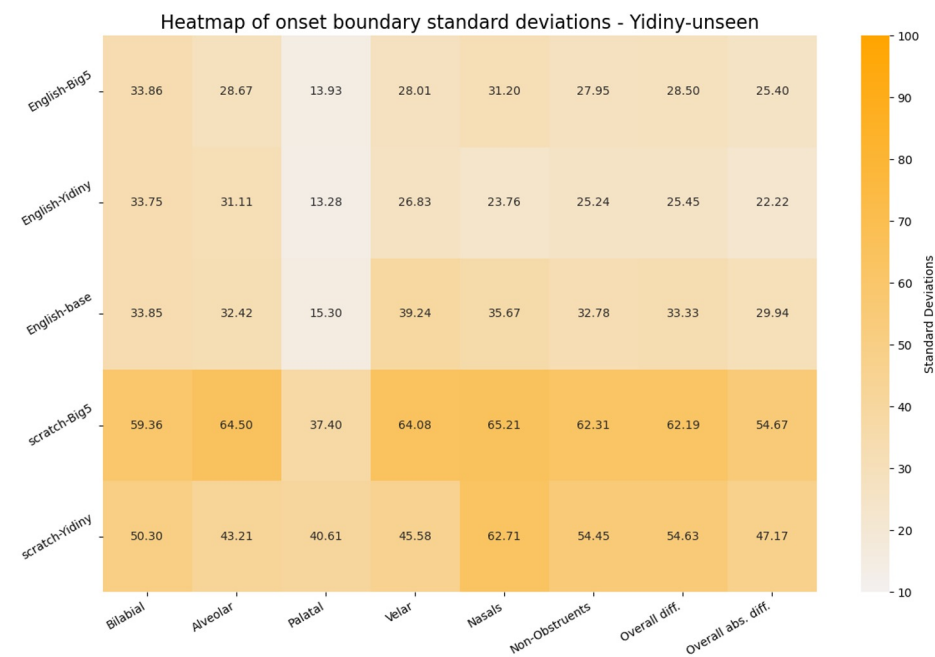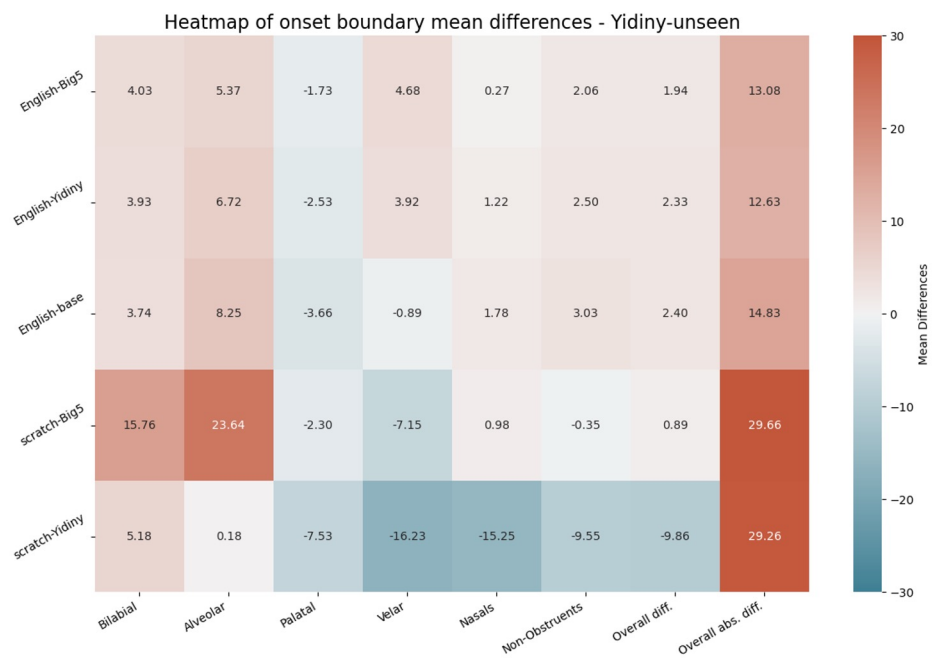Heatmap of onset boundary mean differences - Yidiny-unseen

| | Bilabial | Alveolar | Palatal | Velar | Nasals | Non-Obstruents | Overall diff. | Overall abs. diff. |
|---|---|---|---|---|---|---|---|---|
| English-Big5 | 4.03 | 5.37 | -1.73 | 4.68 | 0.27 | 2.06 | 1.94 | 13.08 |
| English-Yidiny | 3.93 | 6.72 | -2.53 | 3.92 | 1.22 | 2.50 | 2.33 | 12.63 |
| English-base | 3.74 | 8.25 | -3.66 | -0.89 | 1.78 | 3.03 | 2.40 | 14.83 |
| scratch-Big5 | 15.76 | 23.64 | -2.30 | -7.15 | 0.98 | -0.35 | 0.89 | 29.66 |
| scratch-Yidiny | 5.18 | 0.18 | -7.53 | -16.23 | -15.25 | -9.55 | -9.86 | 29.26 |

Heatmap of onset boundary standard deviations - Yidiny-unseen

| | Bilabial | Alveolar | Palatal | Velar | Nasals | Non-Obstruents | Overall diff. | Overall abs. diff. |
|---|---|---|---|---|---|---|---|---|
| English-Big5 | 33.86 | 28.67 | 13.93 | 28.01 | 31.20 | 27.95 | 28.50 | 25.40 |
| English-Yidiny | 33.75 | 31.11 | 13.28 | 26.83 | 23.76 | 25.24 | 25.45 | 22.22 |
| English-base | 33.85 | 32.42 | 15.30 | 39.24 | 35.67 | 32.78 | 33.33 | 29.94 |
| scratch-Big5 | 59.36 | 64.50 | 37.40 | 64.08 | 65.21 | 62.31 | 62.19 | 54.67 |
| scratch-Yidiny | 50.30 | 43.21 | 40.61 | 45.58 | 62.71 | 54.45 | 54.63 | 47.17 |

# Appendix A: PoA Accuracy Kunbarlang (unseen)



Heatmap of onset boundary mean differences - Kunbarlang-unseen

Heatmap of onset boundary standard deviations - Kunbarlang-unseen