Invited talk at International Workshop on Symbolic-Neural Learning

How do audio foundation models understand sound?

Shinnosuke Takamichi (Keio University)



Self introduction



@forthshinji

Name

Shinnosuke Takamichi / 高道 慎之介

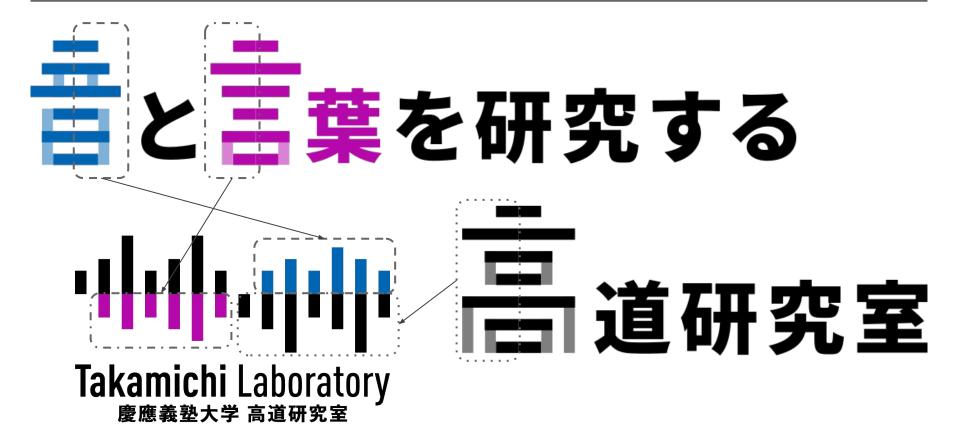
Affiliation

Associate Prof. of Keio University, Japan

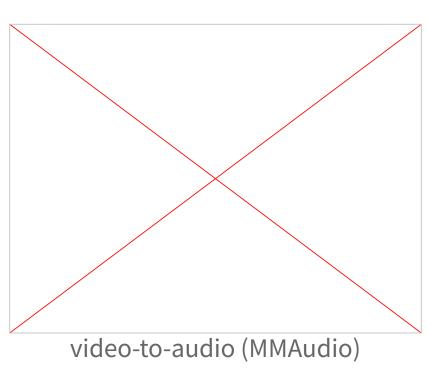
Major

Speech processing

Takamichi (高道) laboratory in Keio, studying audio (音) and language (言葉).



Foundation models for audio (from understanding to synthesis)

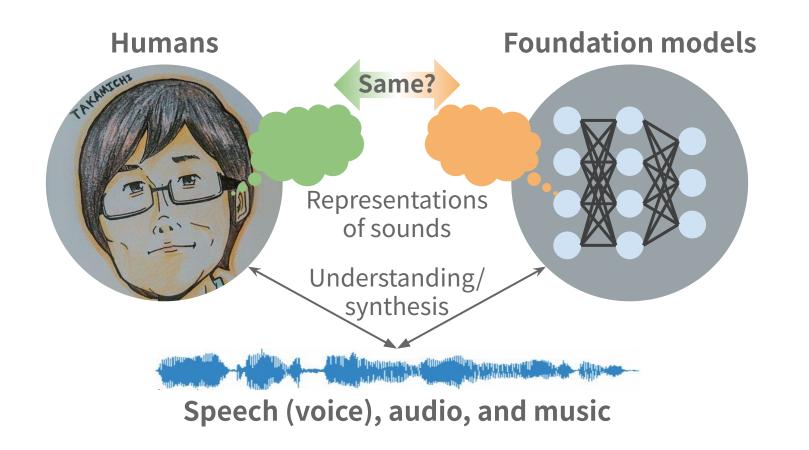






text-to-music (Suno) text-to-speech (E labs)

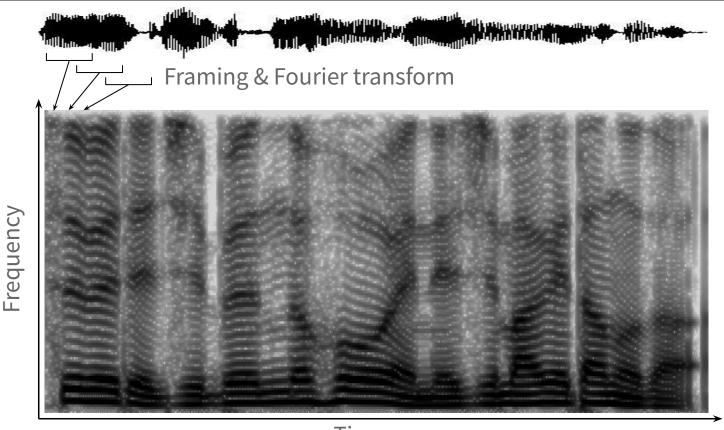
What do the foundation models learn from data?



Topic 1: Language of learned audio symbols

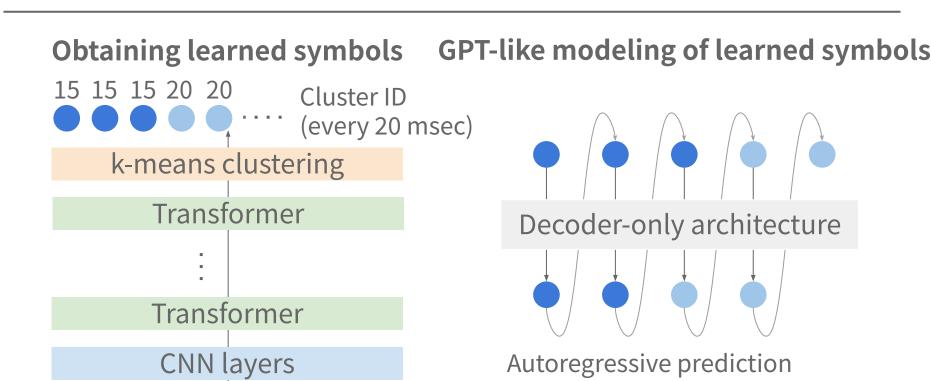
- S. Takamichi et al., "Do learned speech symbols follow Zipf's law?," ICASSP, 2024.
- S. Kando, S. Takamichi et al., "Exploring the effect of segmentation and vocabulary size on speech tokenization for speech language models," Interspeech, 2025.
- J. Park, S. Takamichi et al., "Analyzing the language of neural audio codecs," ASRU, 2025.

Spectrogram: traditional representations of sound



Time

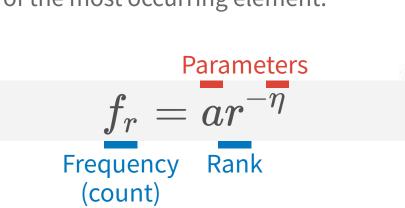
Speech foundation models based on learned audio symbols

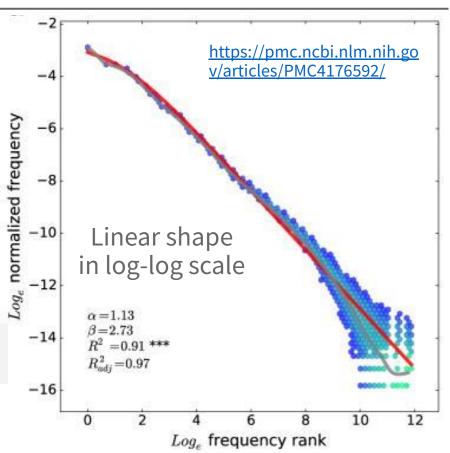


(On a different note) Zipf's law for natural language

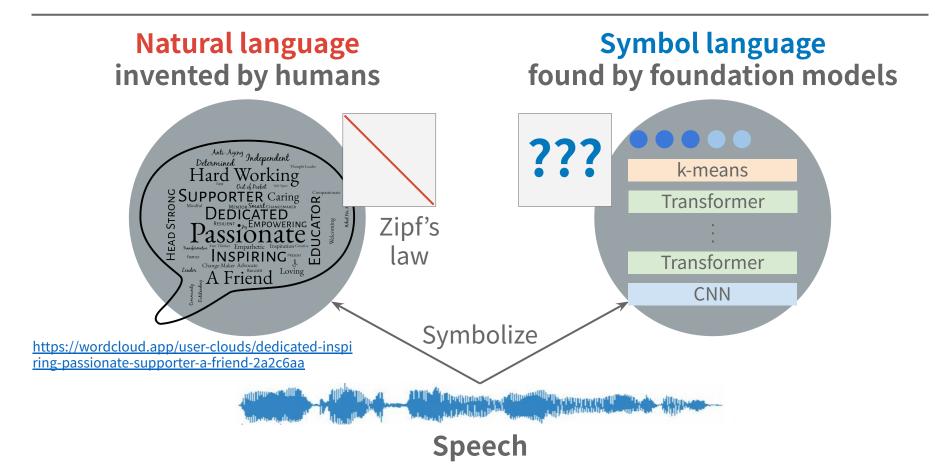
Zipf's law [Zipf49]

- Empirical principle that delineates the frequency of occurrence
- When the occurrence frequency of an element ranks as the *k*-th highest in a set, it equates to 1/*k* of the frequency of the most occurring element.

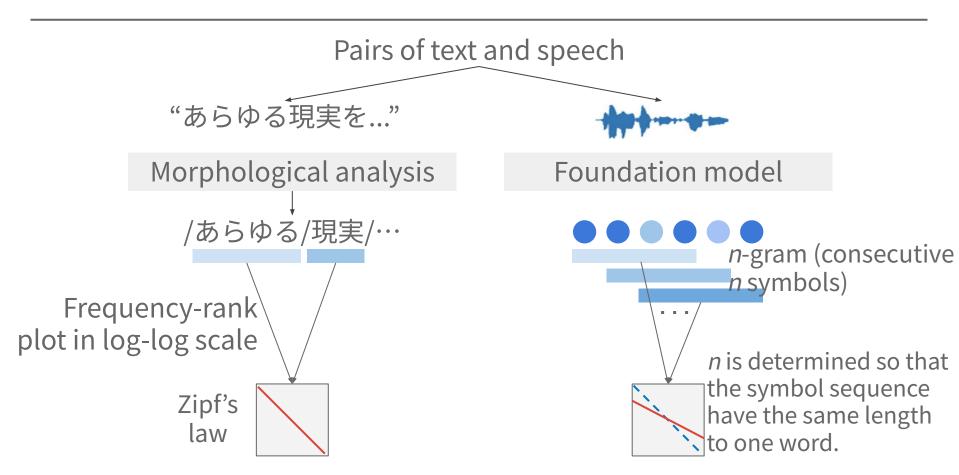




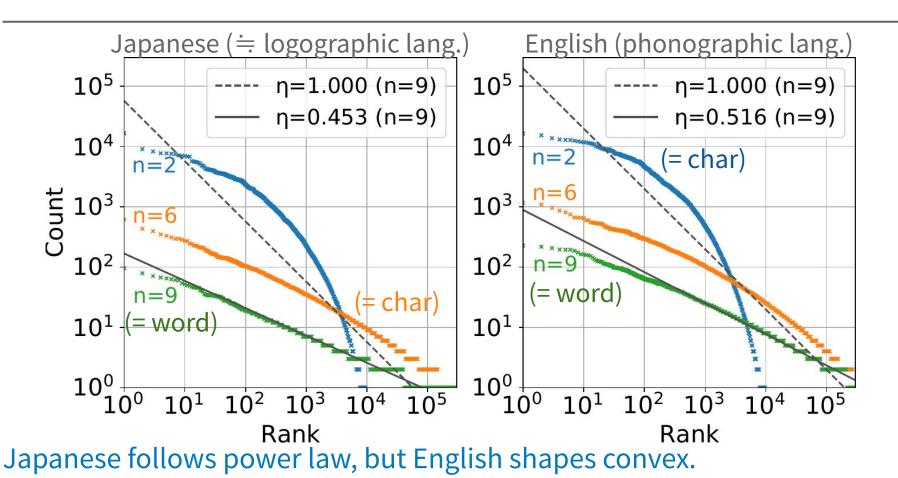
Natural language symbols vs. learned symbols



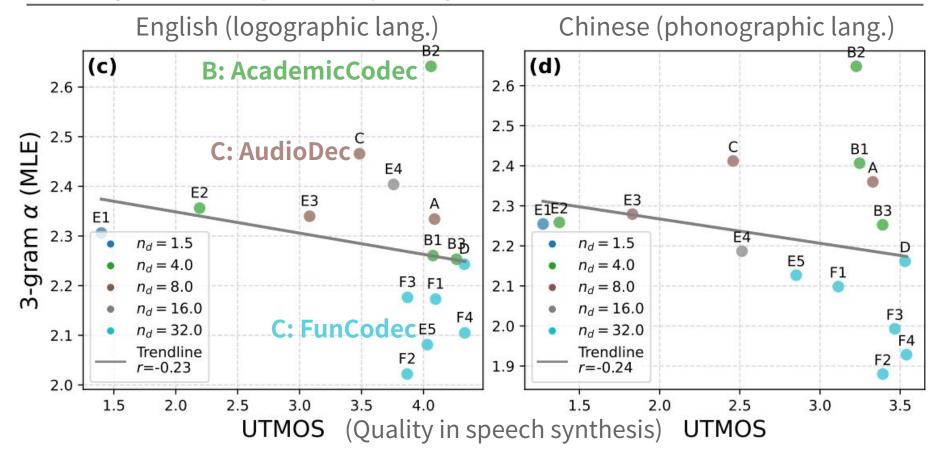
Methodology



Results of learned speech symbols



Parameter ($\alpha = \eta + 1$) in Zipf's law correlates with synthetic speech quality.



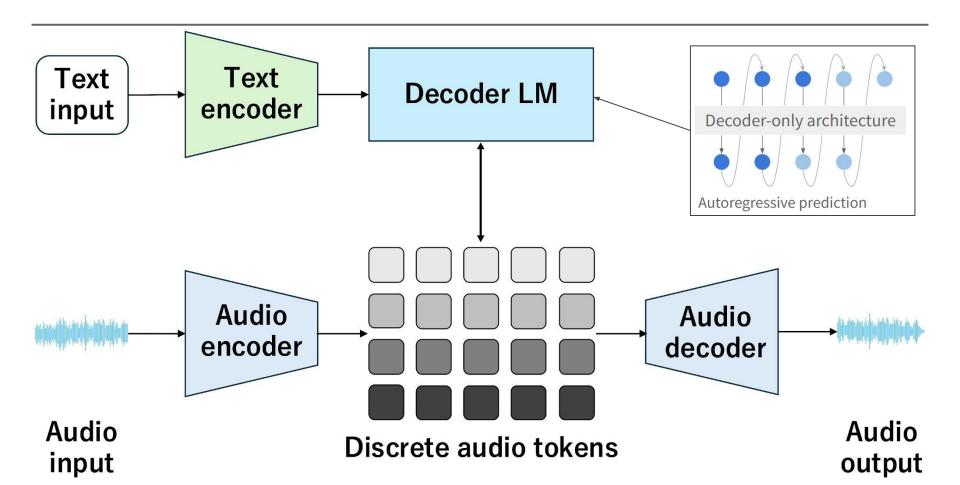
Future work: spoken language is not just audible text, but...

- "Spoken language is not just audible text (= natural language)"
 - Prof. Roger K. Moore (Sheffield Univ.) in Interspeech 2025
 - (not his original phrase)
 - Speech has not only textual information but also tone, pitch, disfluency. timing, etc.
 - Sometimes speech cannot be symbolized natural languages.
- Learned-symbol languages can symbolize any speech and be analyzed by NLP-inspired methods.
 - Zipf's law, Heap's law, etc.

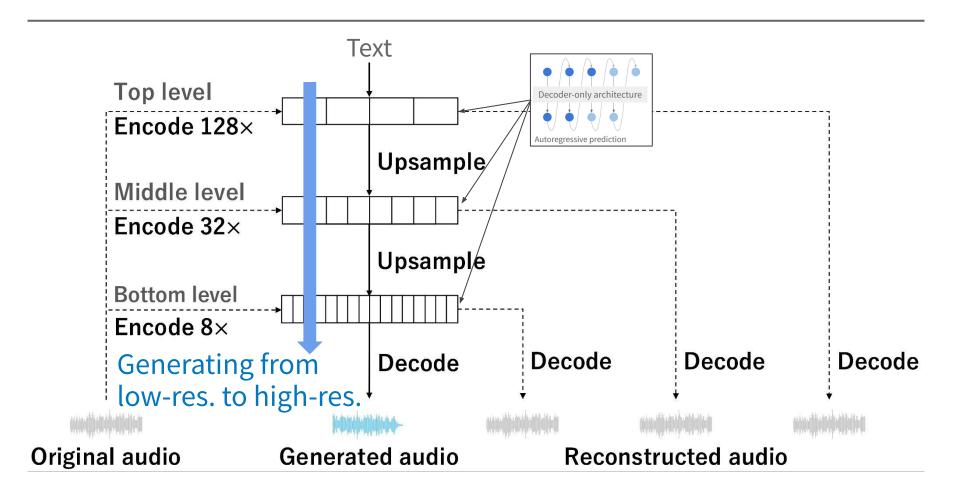
Topic 2: Geometry of music theory in music foundation models

八木 and <u>高道</u>, "音楽基盤モデルは音高情報を螺旋構造に埋め込むか?," MUS, 2025. (as of now, Japanese-language paper only)

Music foundation models: non-hierarchical MusicGen[Copet23]

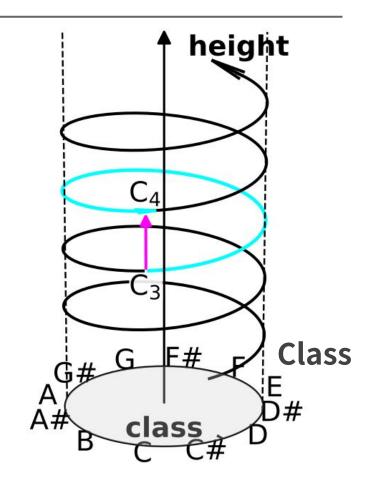


Music foundation models: hierarchical JukeBox[Dhariwal20]

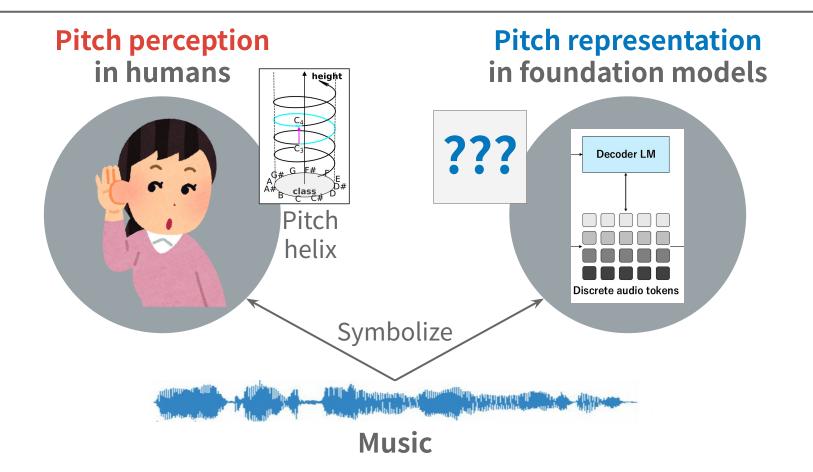


Pitch helix in human perception

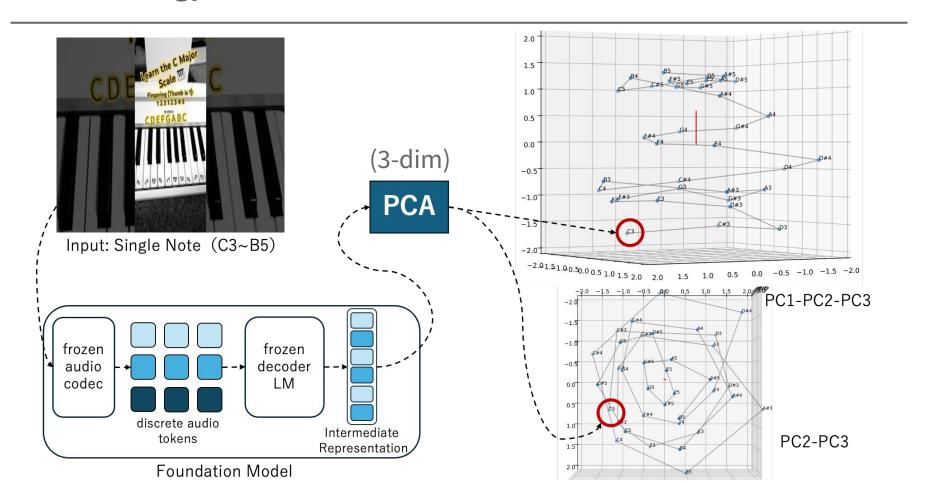




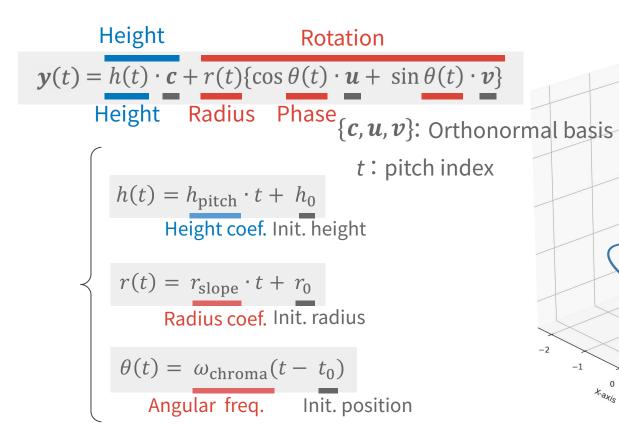
Pitch helix vs. ???

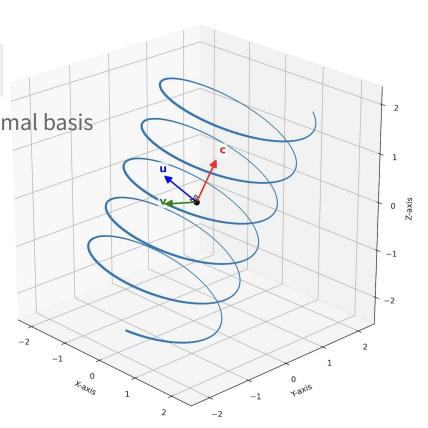


Methodology 1: feature extraction

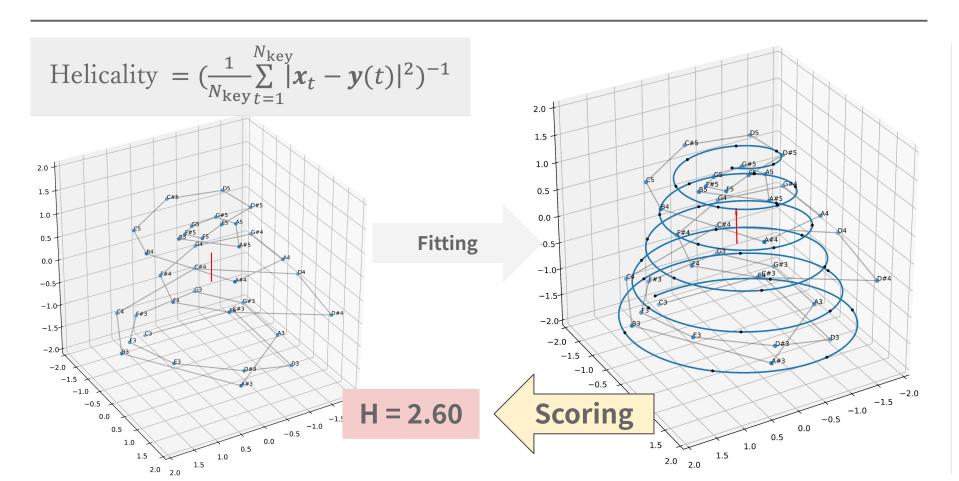


Methodology 2: helix function

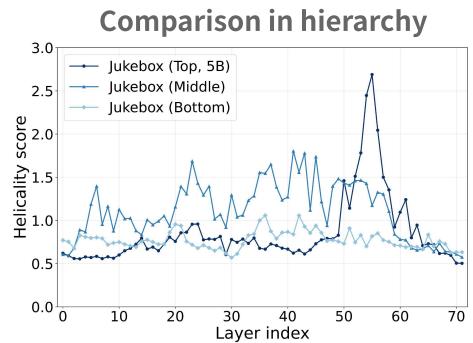




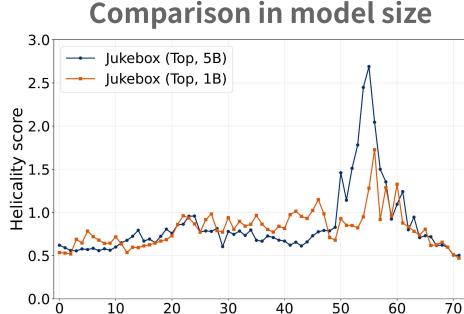
Methodology 3: fitting and scoring



Results in Jukebox (hierarchical model)



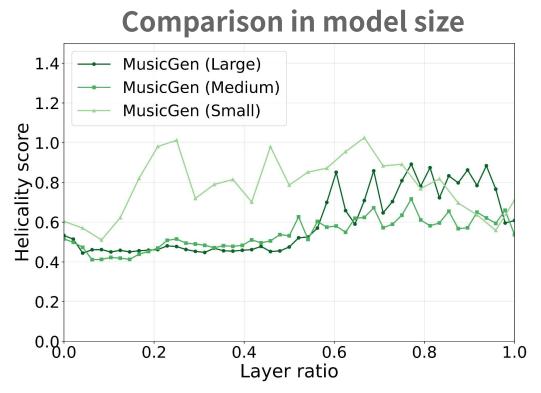
Low temporal model ("Top") has the helix-specific layer.



Larger model (5B) develop layers in which pitch helix becomes clearer.

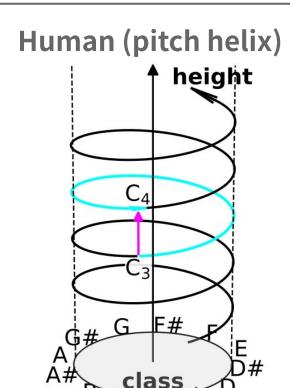
Layer index

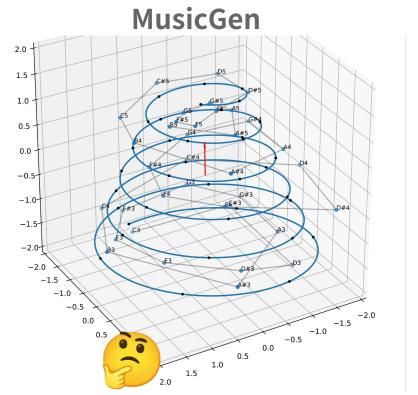
Results in MusicGen (non-hierarchical model)



"Larger model develop layers in which pitch helix becomes clearer" <- observed!

Compassion of perception models





Assuming **one** rotation per octave.

Observing **two** rotations per octave.

Future work

- No more "black box" in music processing.
 - Controllable/explainable music factors.
 - Handling unseen domain data.
 - Too high pitch, speed, tone, color, etc.
- Understanding human perception.
 - Pitch helix really holds true in pitch perception of human?
 - Aligning perception of humans and foundation models.

Conclusion

Conclusion

- Language of learned audio symbols
 - Learned speech symbols instead of hand-crafted features
 - The learned symbols follows the power law.
 - Parameter in the law correlates to speech quality.
- Geometry of music theory in music foundation models
 - Pitch helix, a pitch perception model for human.
 - Fitting helix-inspired function to data in foundation models.
 - The foundation models also follows the pitch helix.