

# Measuring Sound Symbolism in Audio-visual Models

Wei-Cheng Tseng\*, Yi-Jen Shih\*, David Harwath and Raymond Mooney

Department of Computer Science, The University of Texas at Austin, USA

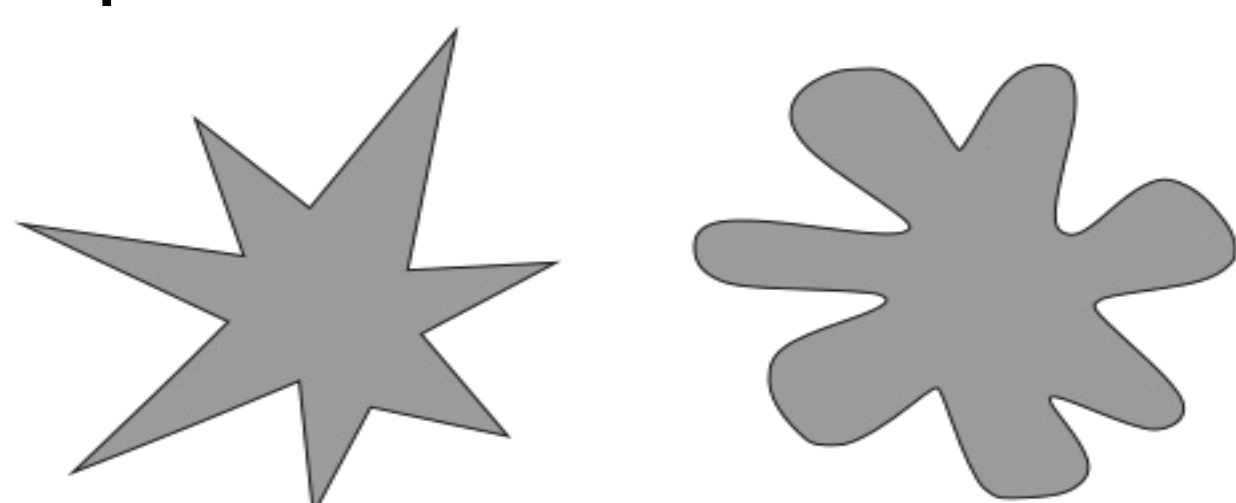
## TL;DR

This work aims to bridge the gap between cognitive science and speech technologies. We investigate the existence of a linguistic phenomenon—**sound symbolism**—in pre-trained audio-visual models. Experimental results show that some models can capture sound-meaning associations akin to those observed in humans, influenced by the pre-training methods and datasets.

## Background & Motivation

- Audio-visual representation learning has become an important research topic
- Pre-trained audio-visual models emerges intriguing language processing capabilities without direct supervision:
  - AV-HuBERT [1] → better acoustic unit discovery
  - VG-HuBERT [2] → word discovery & segmentation
  - AV-NSL [3] → syntax acquisition
- **What about more abstract linguistic phenomena?**
- We look at sound symbolism. Specifically, **“Kiki-bouba Effect”**
  - Independent of cultural and linguistic background
  - Important for early language acquisition

## Dataset Collection



Audio:

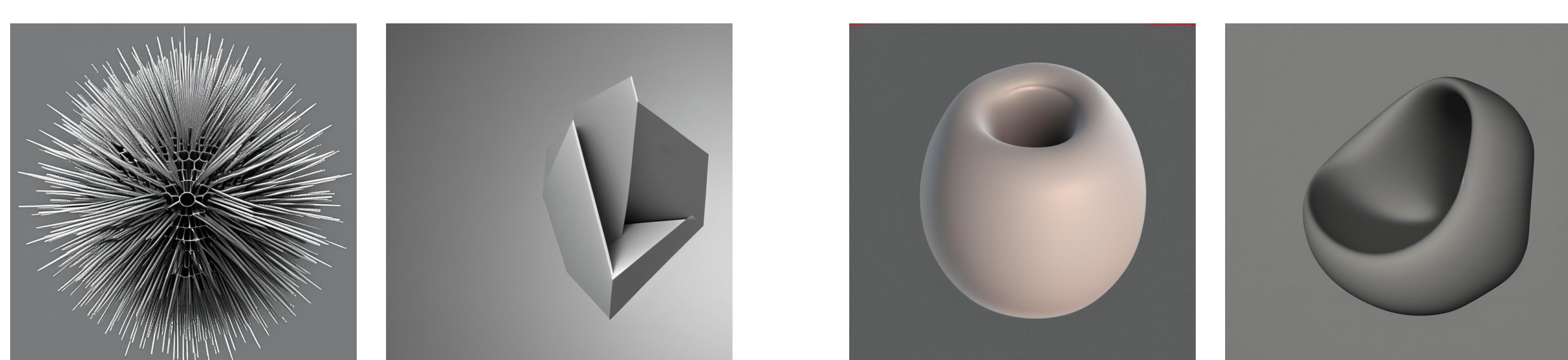
- Synthesize sounds from predefined sets [4] of **sharp** or **round** phones w/ certain template:

$$\begin{aligned}
 C_{\text{SHARP}} &= \{k t p f d_3 z\} & V_{\text{SHARP}} &= \{\epsilon i:\} \\
 C_{\text{ROUND}} &= \{m n l b d g\} & V_{\text{ROUND}} &= \{o: u:\} \\
 C_{\text{NEUTRAL}} &= \{f s v\} & V_{\text{NEUTRAL}} &= \{a:\}
 \end{aligned}$$

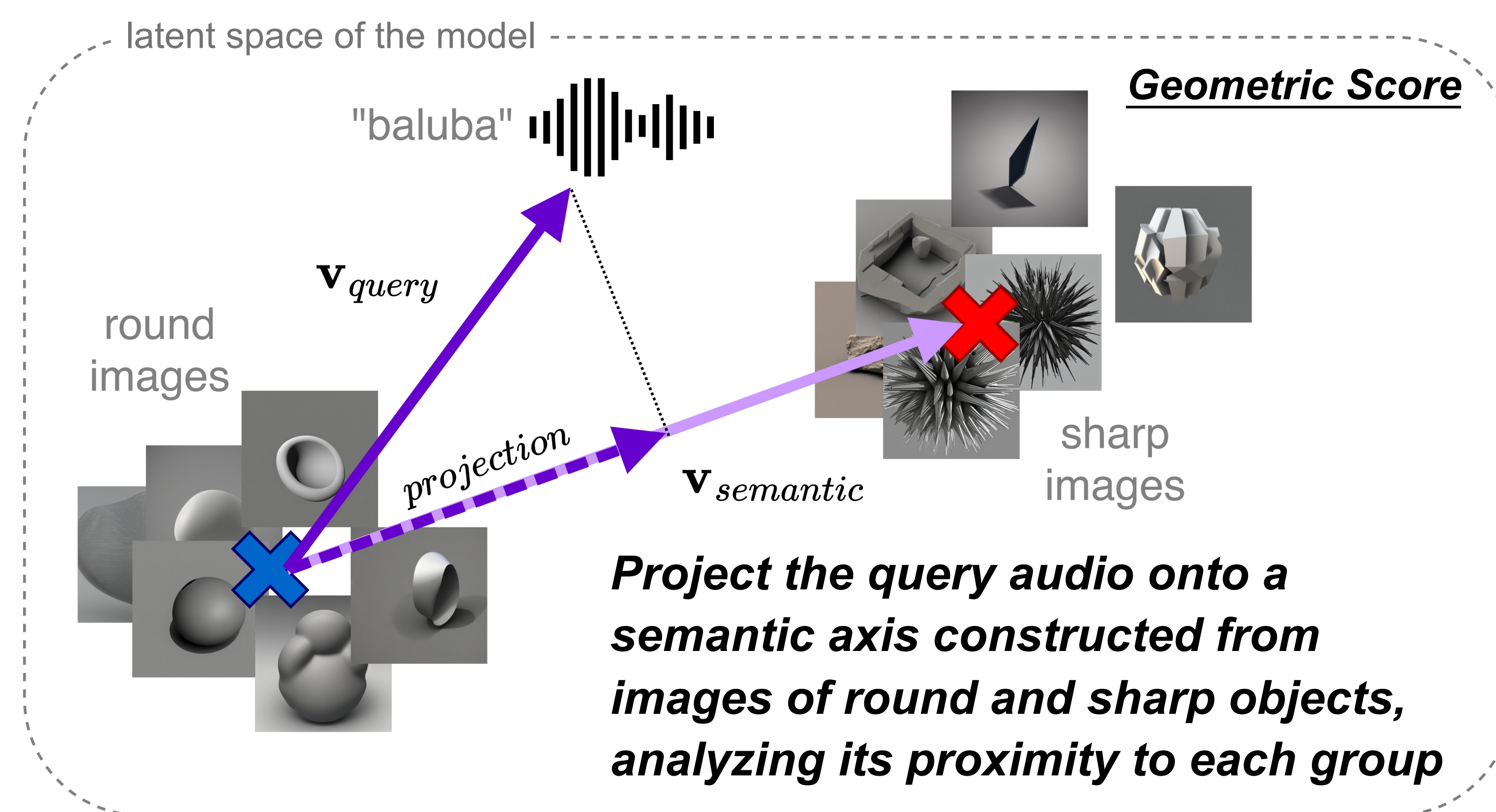
- $(C_1)(C_2)(C_1)$
- Can't draw from opposite groups at same time

Image:

- Define sets of **sharp** or **round** adjectives:
  - $W_{\text{SHARP}} = \{\text{sharp spiky angular ... pointed rugged}\}$
  - $W_{\text{ROUND}} = \{\text{round circular soft ... smooth chubby}\}$
- Synthesize w/ certain template
- **A 3D-rendering of a  $\langle w \rangle$  object**

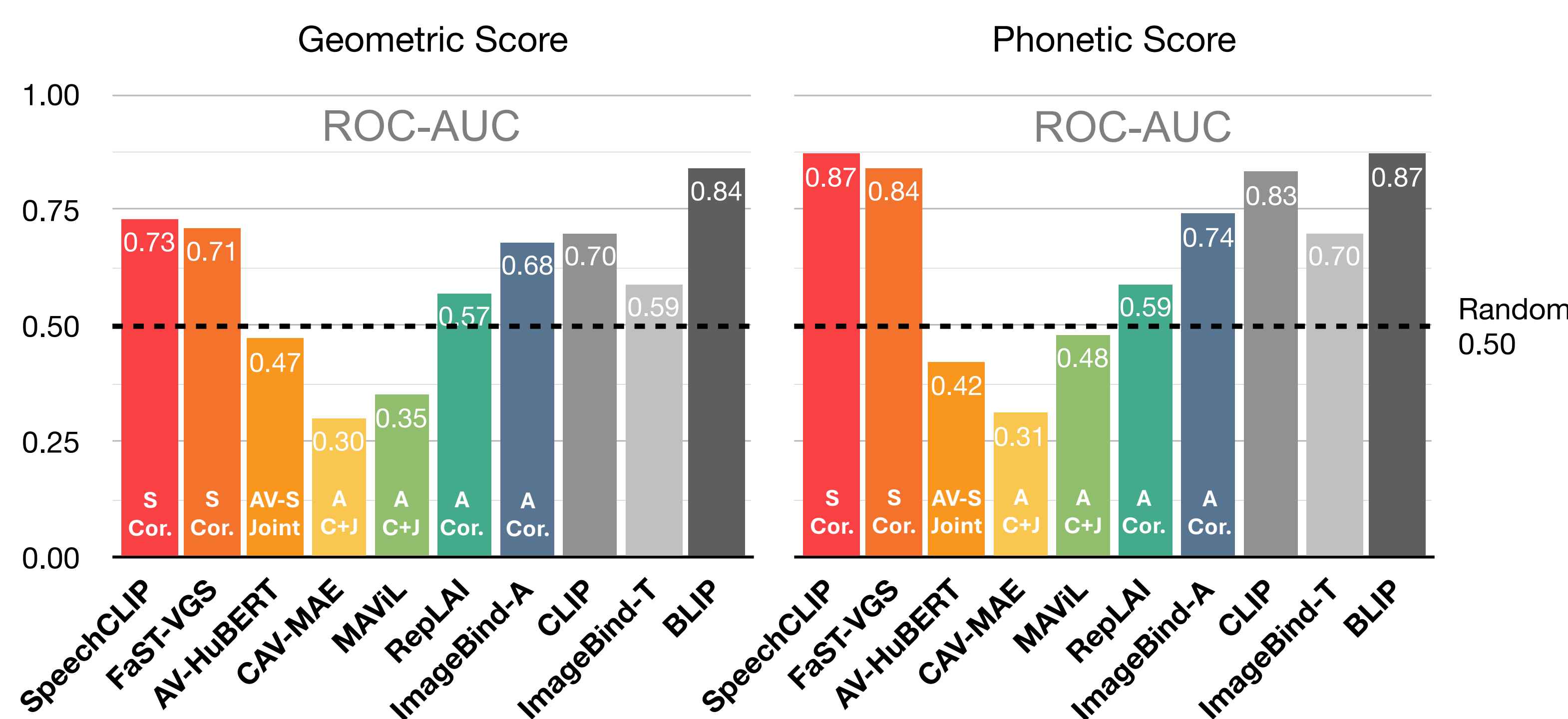


## Evaluation: Zero-shot Probing

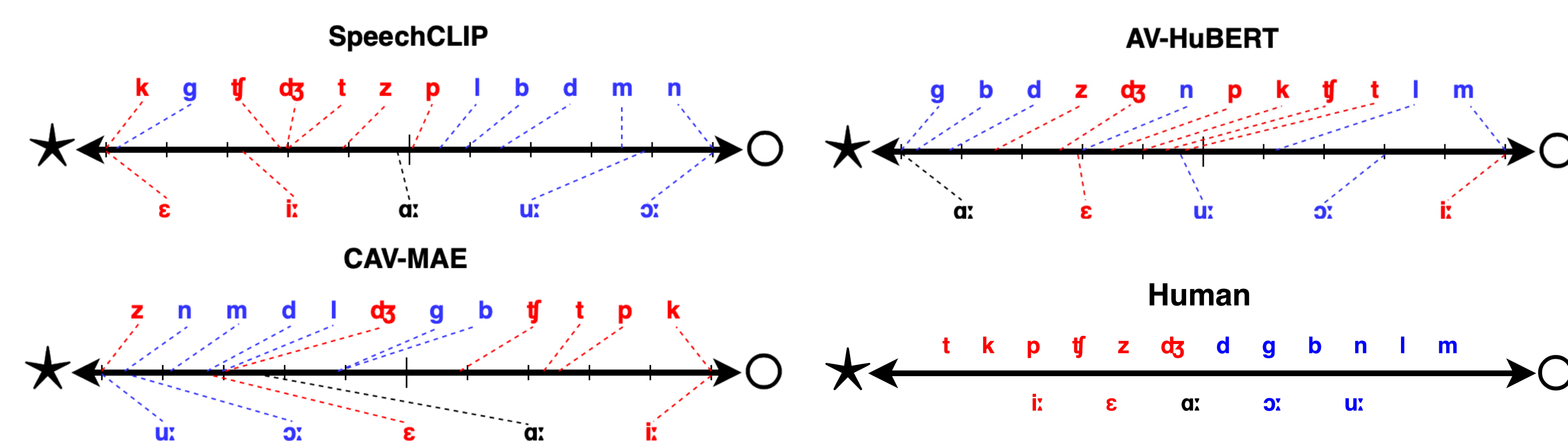


## Results

- Apply a threshold to the scores (magnitude of the projected vectors) for binary classification



- Models pretrained on **spoken image captions** show more profound association than those trained on **general audio**
- Pretraining objectives/methodologies also affect
- Ranking of proximity of phone to sharp/round attr.



## Contributions & Conclusions

- Investigate presence of **kiki-bouba effect** in pre-trained audio-visual models
- Reveal significant corr. between patterns of sound symbolism and models learned on spoken captions, analogous to human learning process
- Support non-arbitrariness of language and provide insight into audio-visual learning algorithm

## Reference

- [1] B. Shi et al., "Learning audio-visual speech representation by masked multimodal cluster prediction," 2022.
- [2] P. Peng et al., "Word Discovery in Visually Grounded, Self-Supervised Speech Models," 2023.
- [3] C.J. Lai et al., "Audio-visual neural syntax acquisition," 2023.
- [4] Kelly McCormick et al., "Sound to meaning mappings in the bouba-kiki effect.," 2015.

Paper link

