

Speech models represent
phonemes in an emergent,
substance-informed feature space

INTRODUCTION

To what degree do featural representations reference phonetic substance?

Substance-informed (Chomsky & Halle 1968)

Substance-free (Hale & Reiss 2000)

Phonemes contain information about their phonetic characteristics

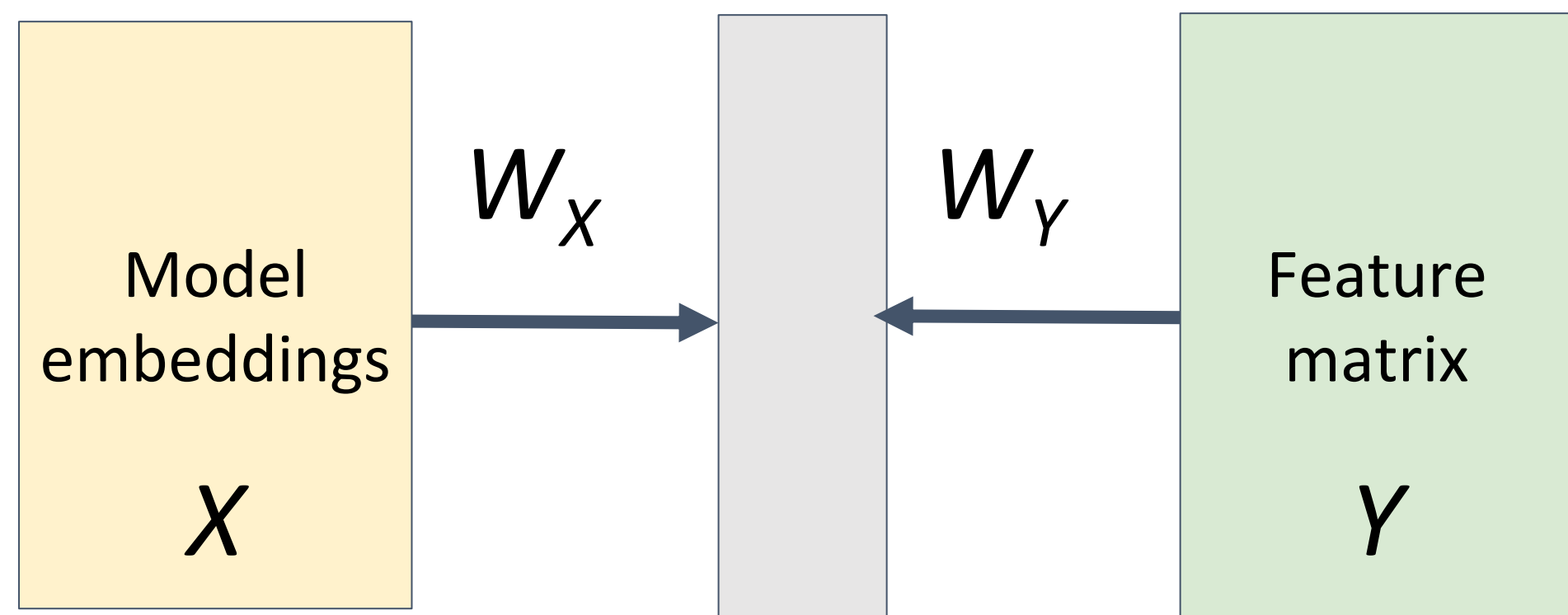
Phonemes are represented distributionally, not related to their phonetic characteristics

/ŋ/, /g/, /m/ are phonetically similar

/ŋ/, /ʒ/, /v/ are distributed similarly

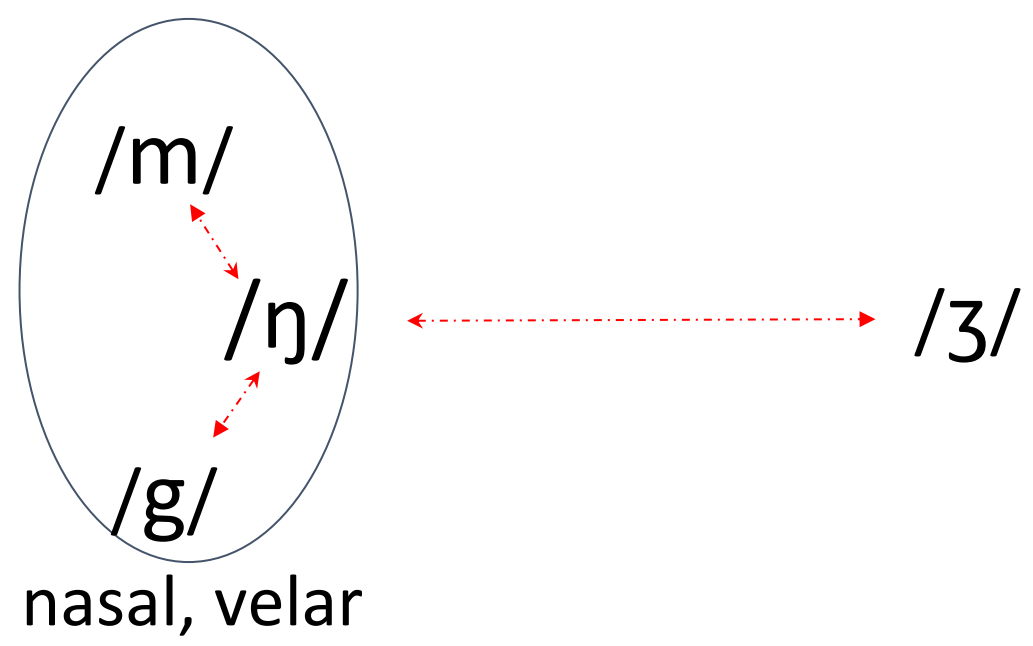
Which type of feature system better corresponds to models' representations?

APPROACH

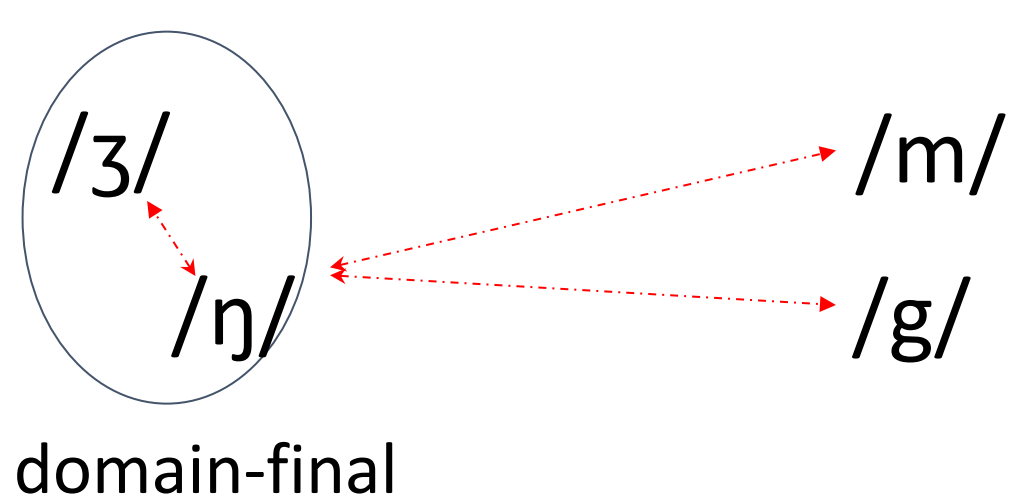


Align Speech Model **embeddings** of phonemes with different **feature systems** using Canonical Correlation Analysis (CCA):

- **Substance-informed feature system from Hayes (2009) (17 dim.):** → what do phonemes sound like?



- **Substance-free feature system** from Mayer (2020) (34 dim.): → how are phonemes distributed?



Data

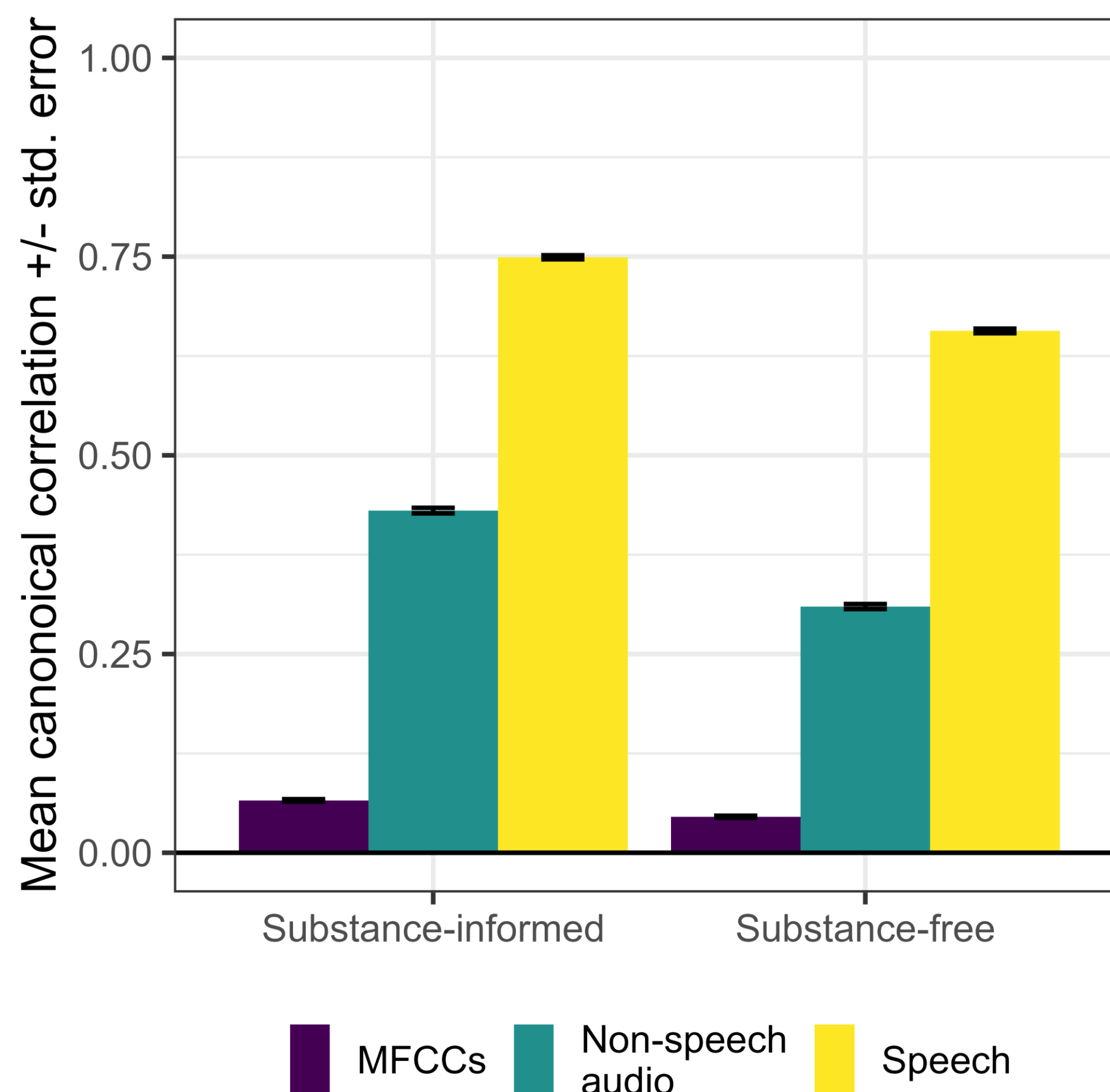
39 English phonemes, extracted from CV and VC sequences of English, synthesized by 10 TTS voices.

MODELS TESTED

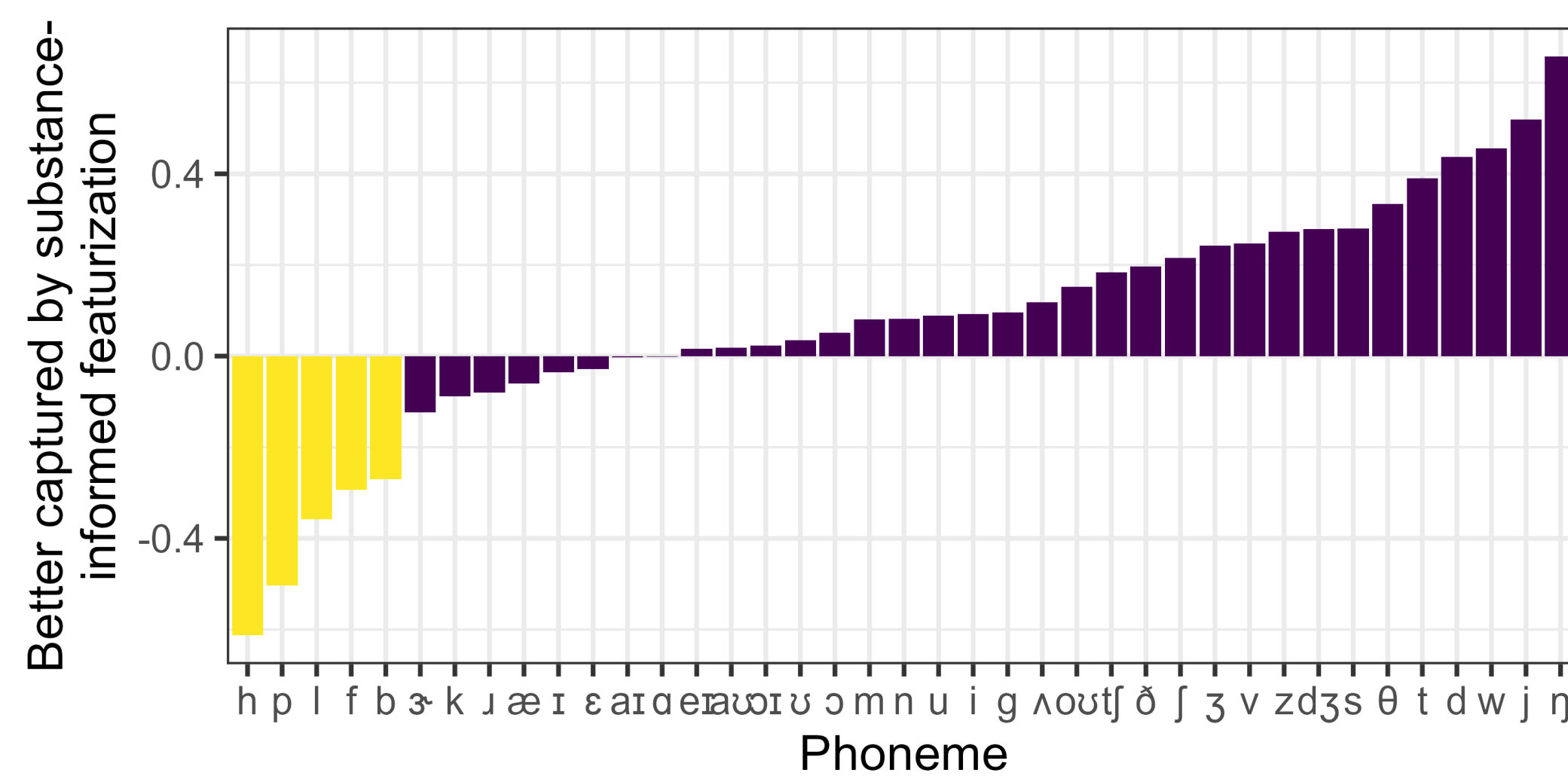
Hubert-Large (trained on English);
Hubert-Large (trained on non-speech
ambient sounds); Wav2Vec2-Large (not
shown)

RESULTS

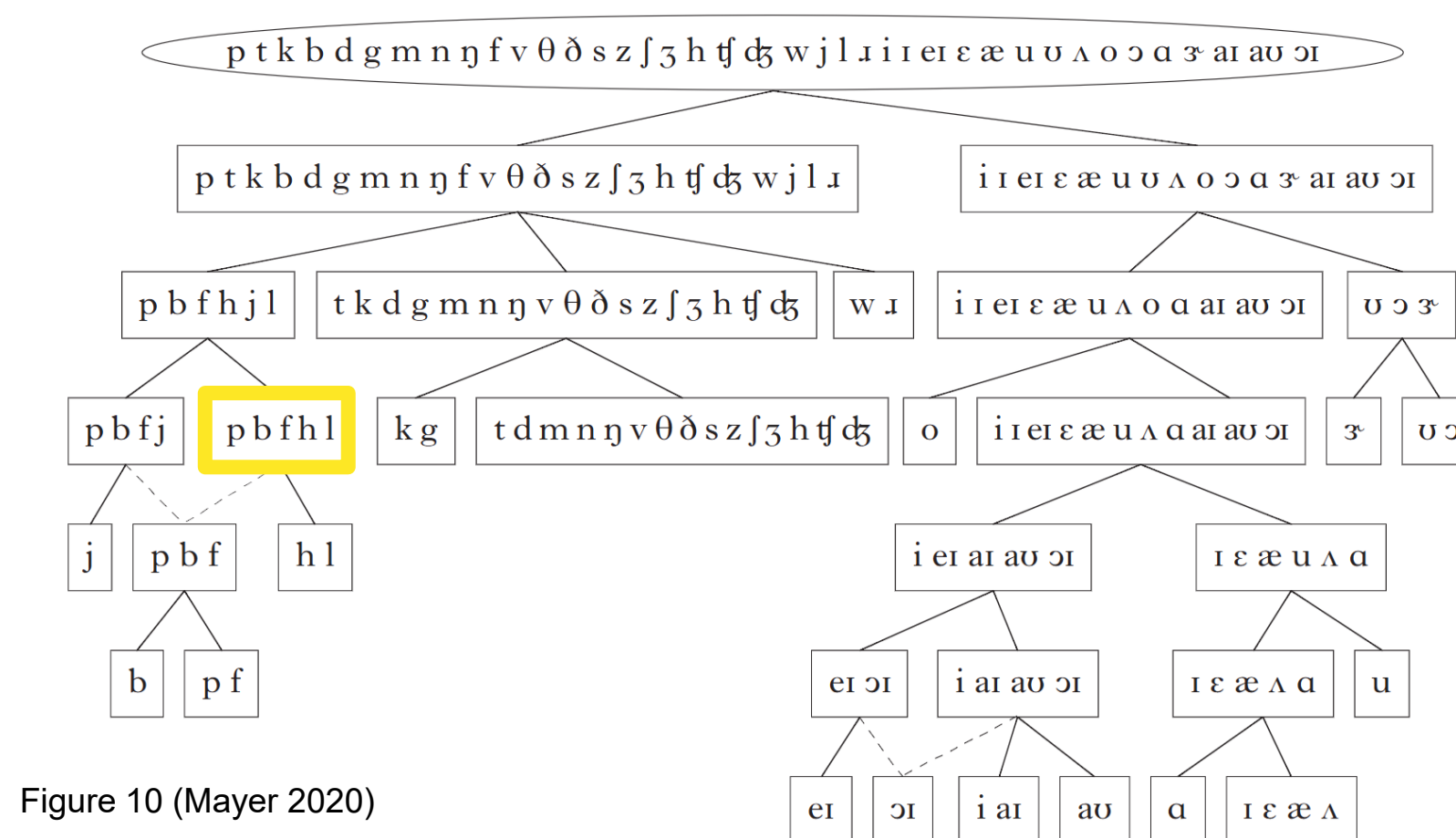
CCA model-feature alignment:



Second-order representational similarity:



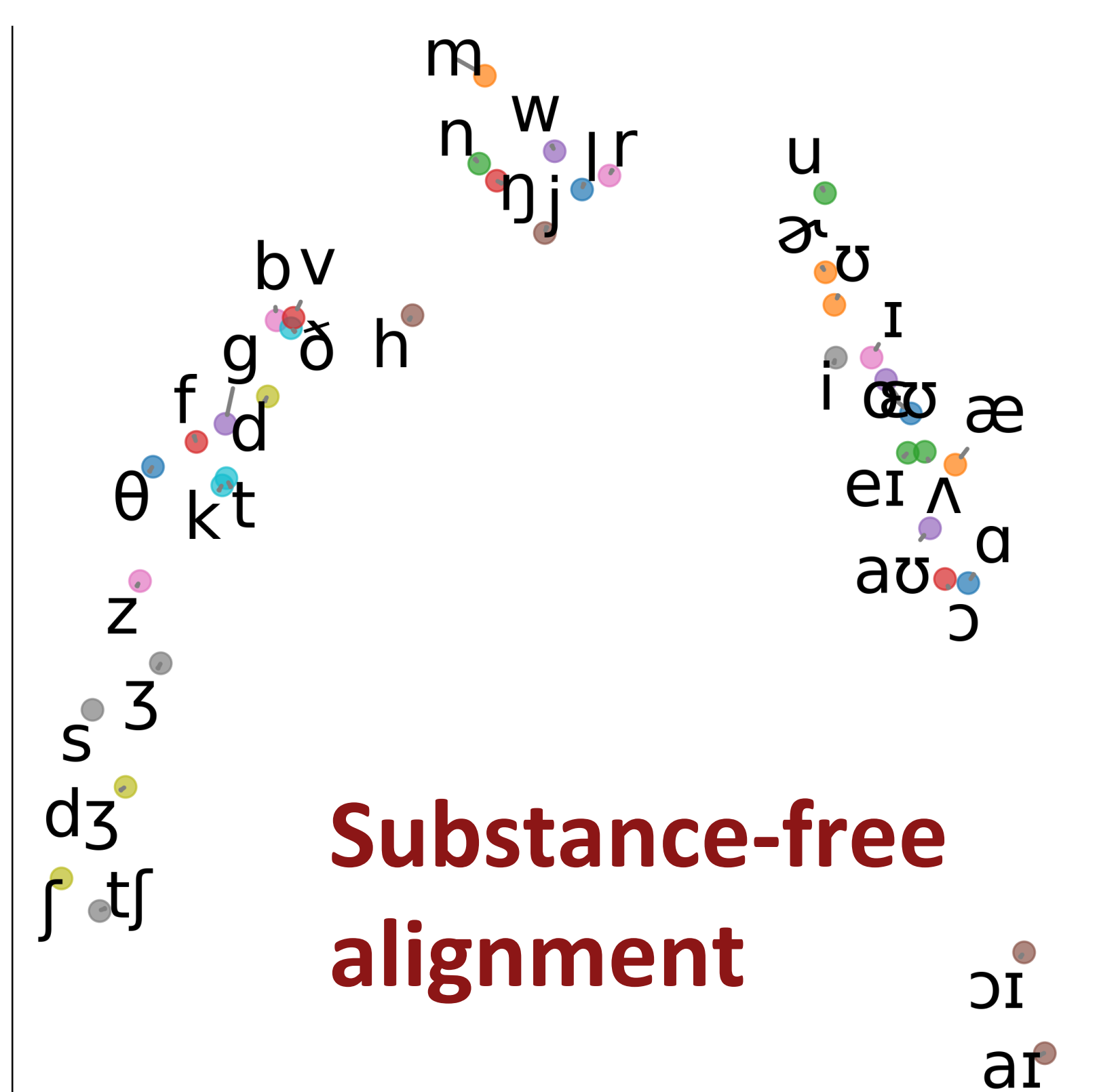
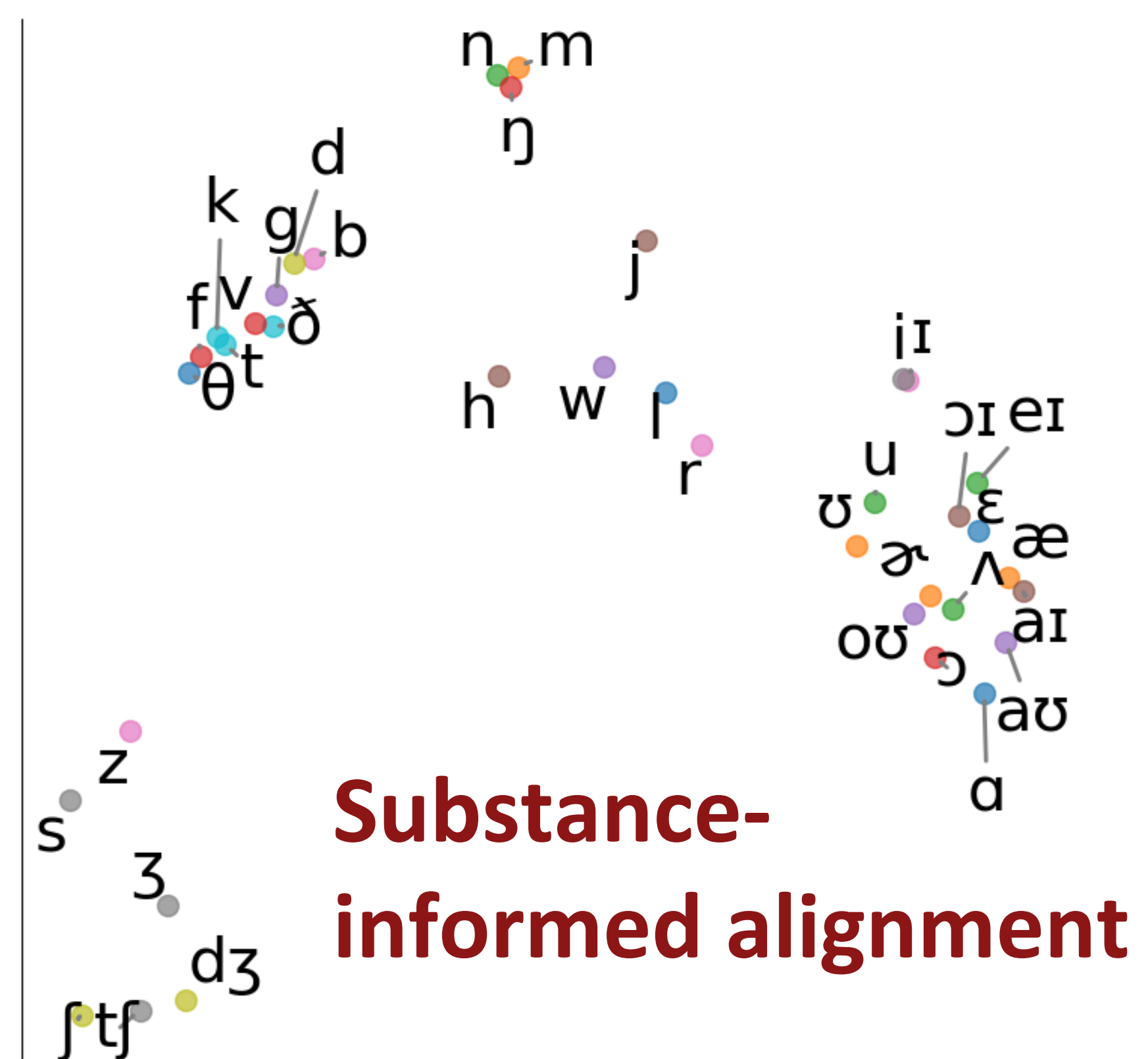
Highly correlated with preference for onset position in Mayer (2020)



12th (for substance-informed features)
and 14th (substance-free) layers shown.

ERROR ANALYSIS

Model embeddings struggle to capture the [η] ~ [ɜ] relationship in the substance-free feature system, even with supervision through CCA.



TAKEAWAYS

Speech models exhibit **emergent, featurally-structured** phoneme representations

Speech model's phoneme embeddings are primarily **substantive**: they encode phonetic properties better, on average, than abstract distributional properties.

