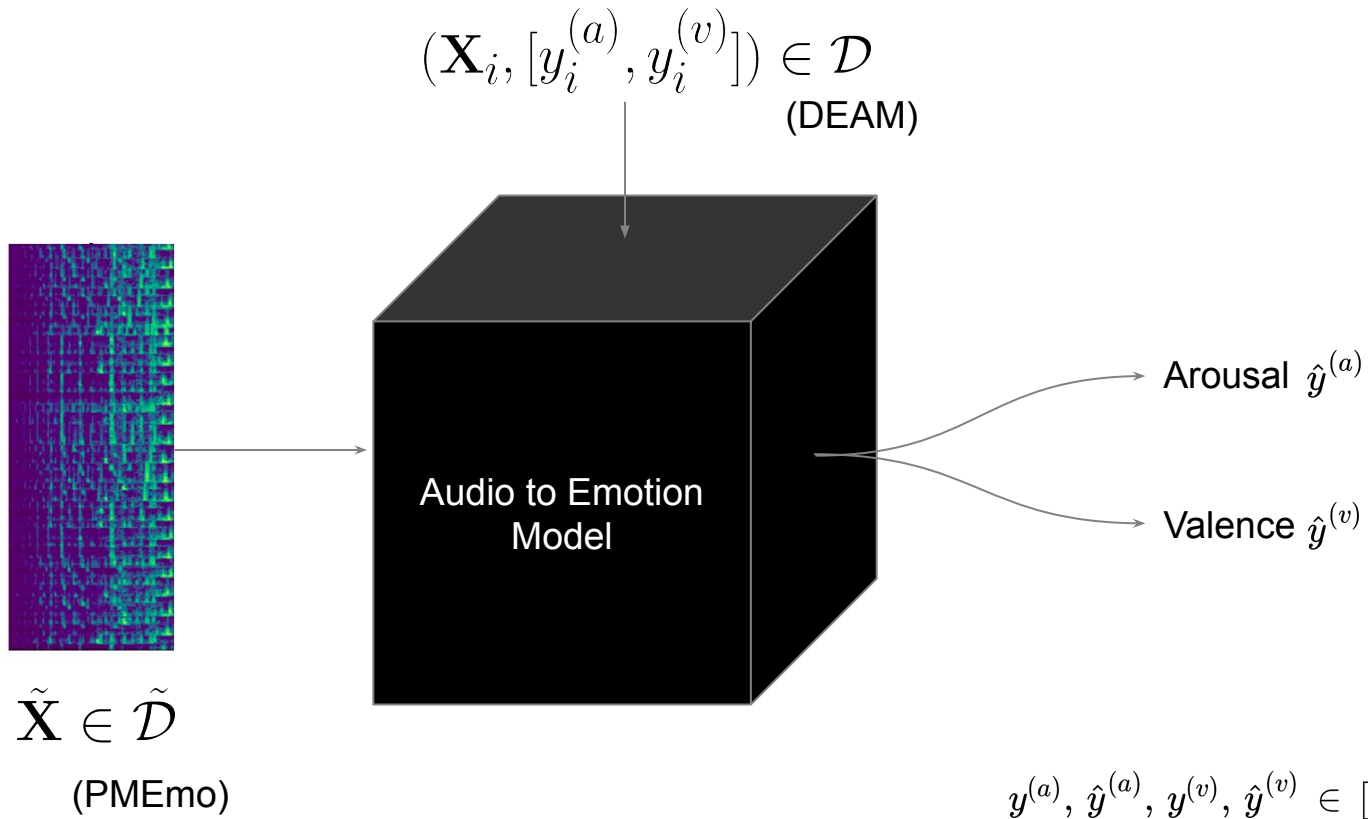


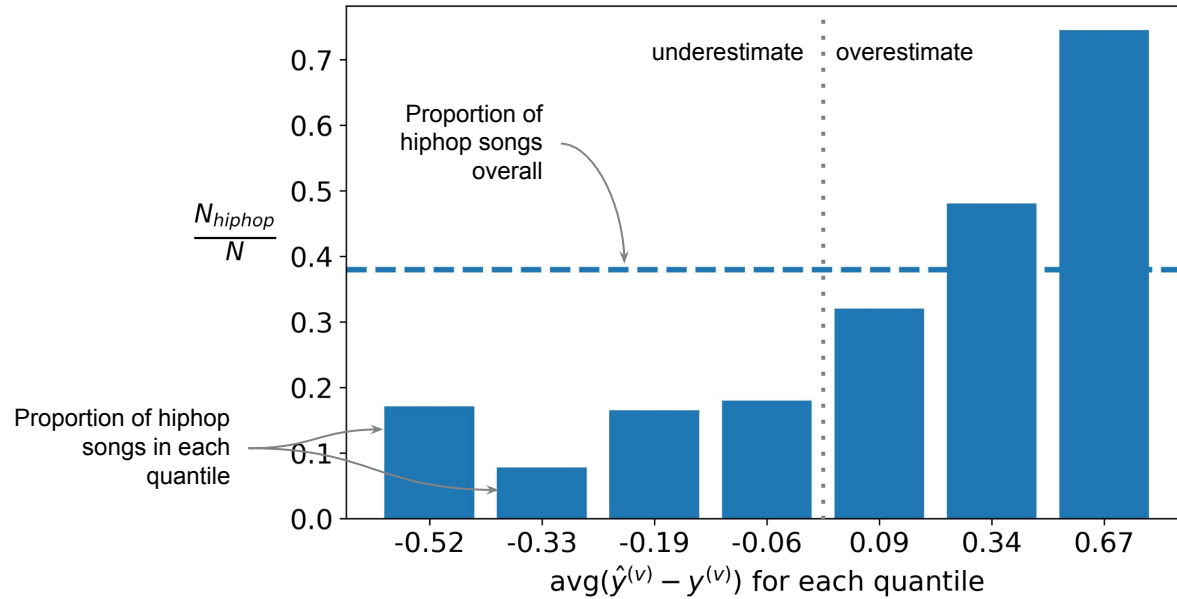


Tracing Back Music Emotion Predictions To Sound Sources And Intuitive Perceptual Qualities

Shreyan Chowdhury, Verena Praher, Gerhard Widmer



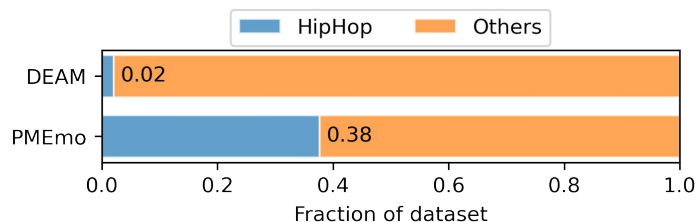
Valence Error Quantiles



Test Set: PMEmo

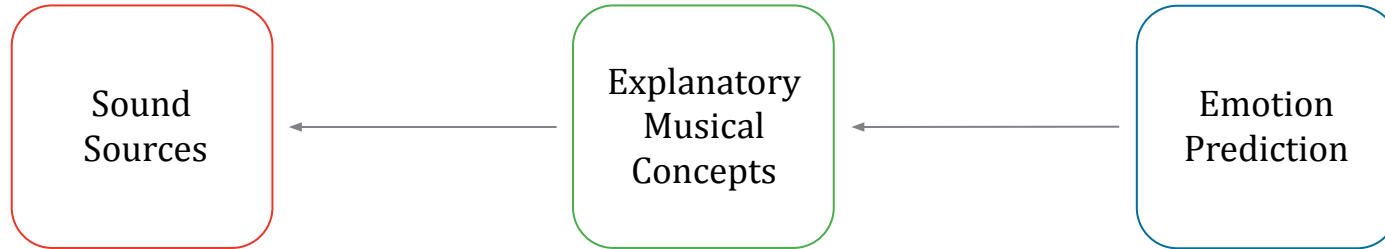
Peering Inside the Model and the Music

- Why the overestimation?

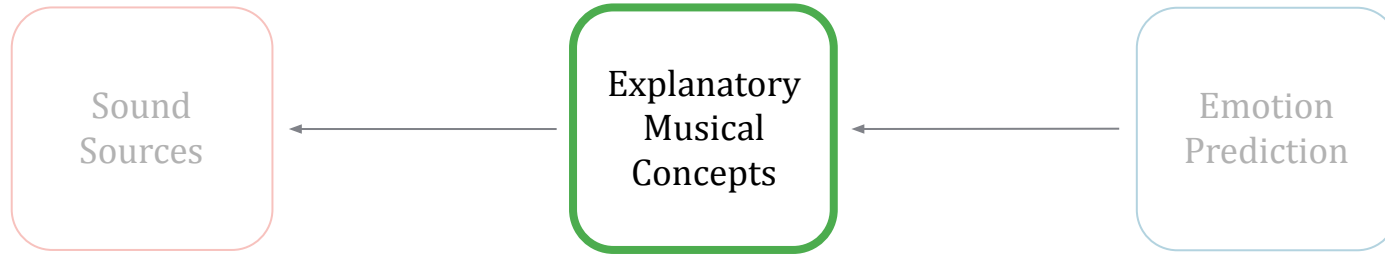


- What musical properties lead to this specific bias?
 - Two-level explanations
- Does intervention lead to expected change?
 - Retraining on balanced data

Two Levels of Explanations



Two Levels of Explanations



Explanatory Musical Concepts: Mid-level Features

Low-level features,
such as pitch

Building blocks of musical
signals

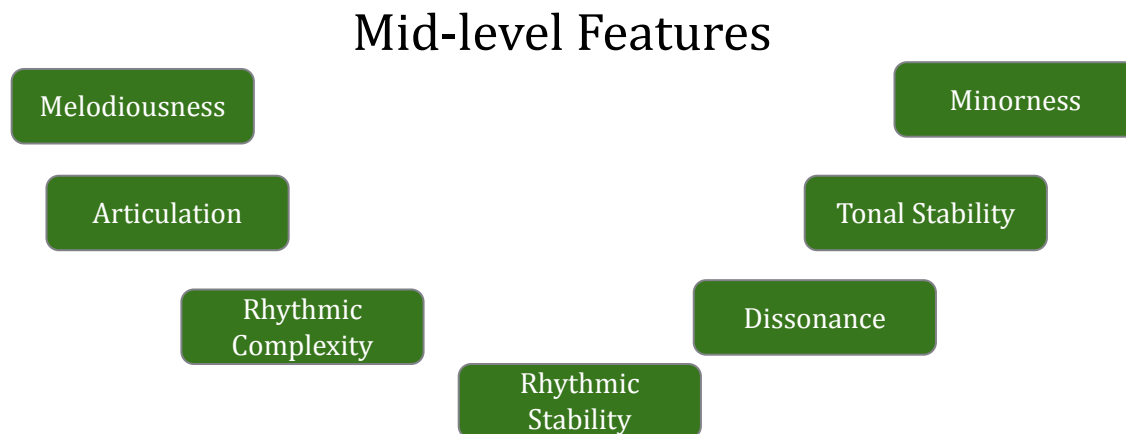
Mid-level Features

Perceptual and subjective, but
make intuitive musical sense

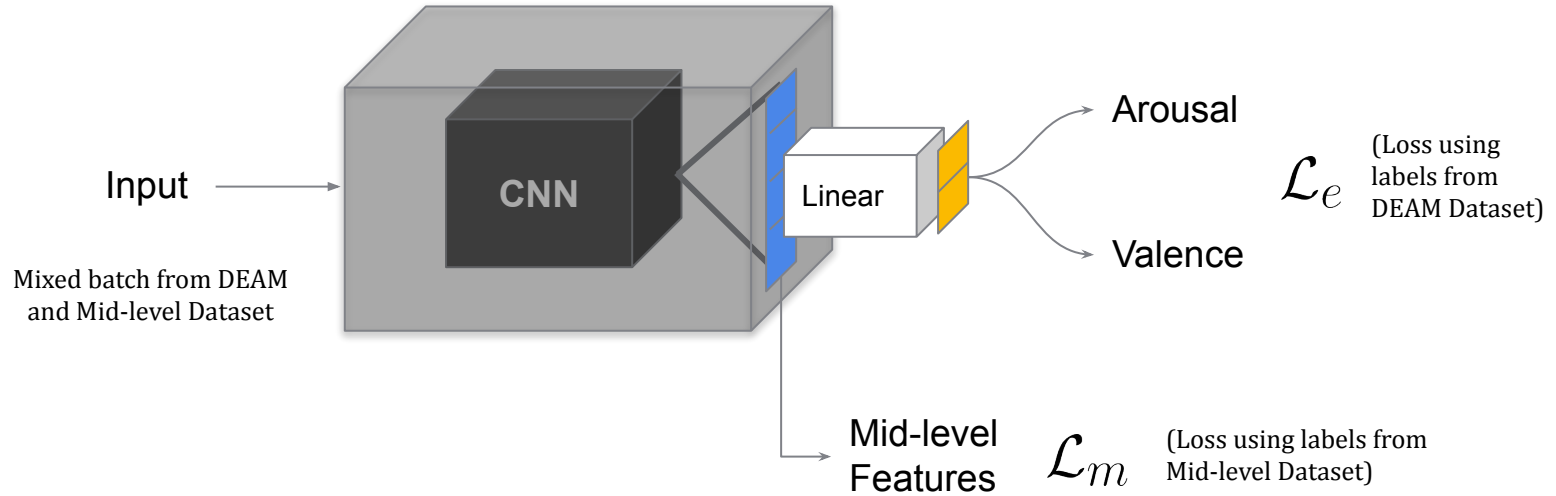
High-level
features, such as
genre

Subjective, abstract
descriptions

Mid-level Features: What are they?



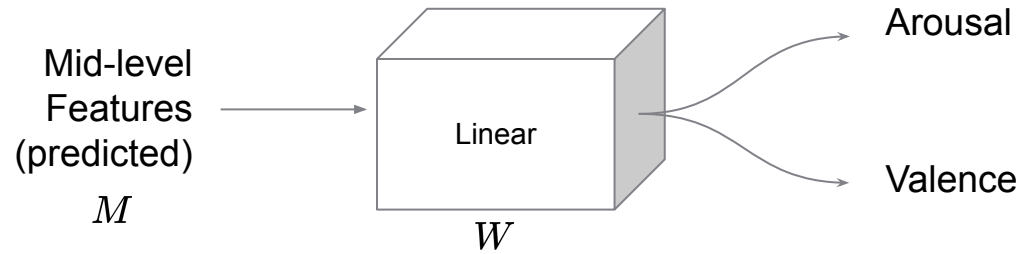
Training the Explainable Model



Joint multitask learning: backpropagate combined loss

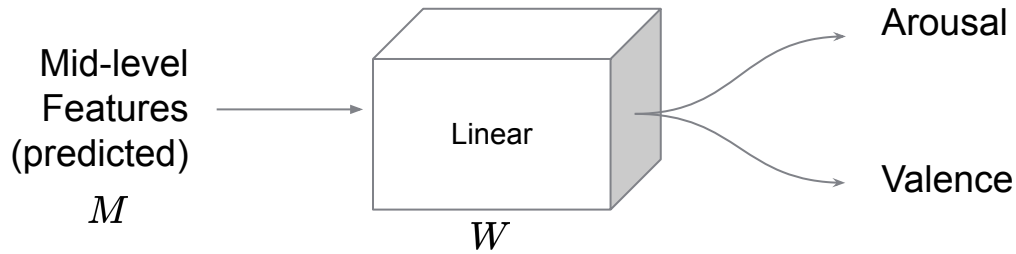
$$\mathcal{L} = \mathcal{L}_m + \mathcal{L}_e$$

Obtaining Mid-level Explanations



Effect:
$$E_j^{(k)} = w_j^{(k)} m_j \quad k \in \{\text{arousal, valence}\}$$

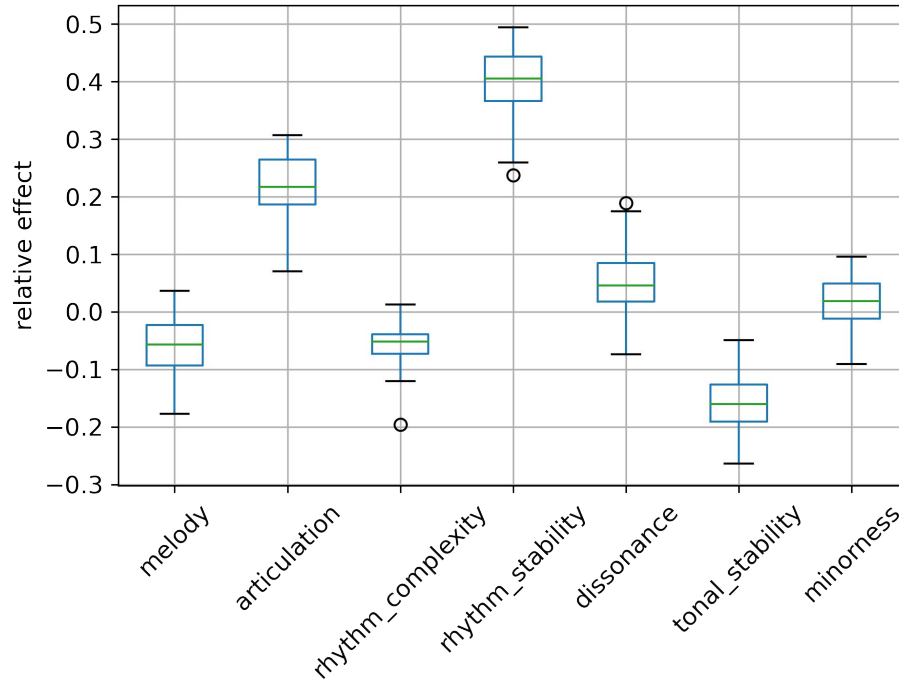
Obtaining Mid-level Explanations



Effect:
$$E_j^{(k)} = w_j^{(k)} m_j \quad k \in \{\text{arousal, valence}\}$$

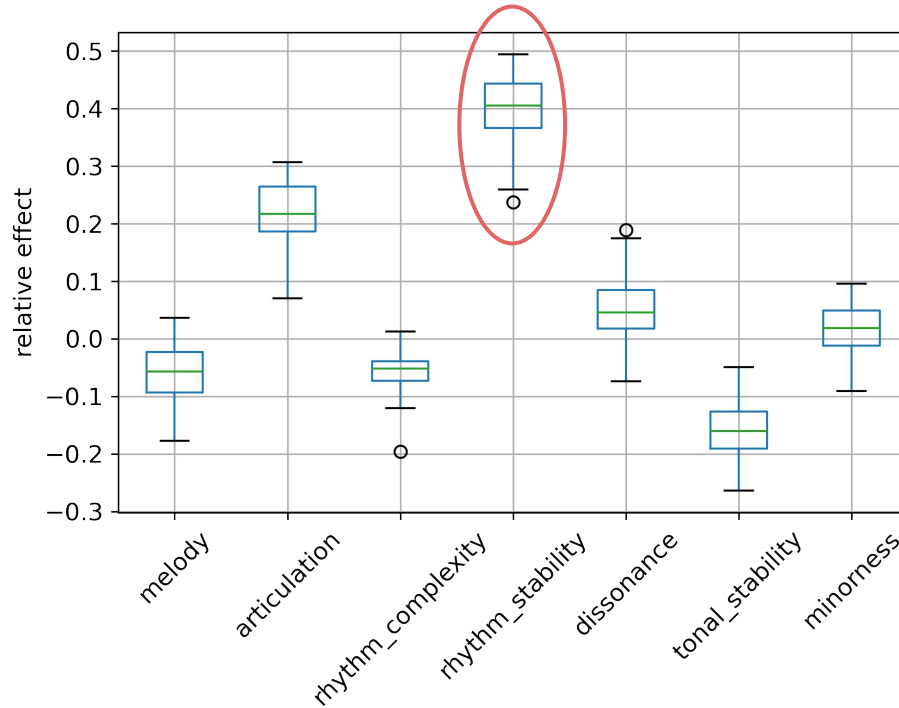
Relative Effect:
$$E_{\text{rel } j}^{(k)} = \frac{|E_j^{(k)}|}{\sum_i |E_i^{(k)}|}$$

Mid-level Explanations for Valence in Hip-hop



Relative effects of
mid-level features on
valence
for hip-hop songs

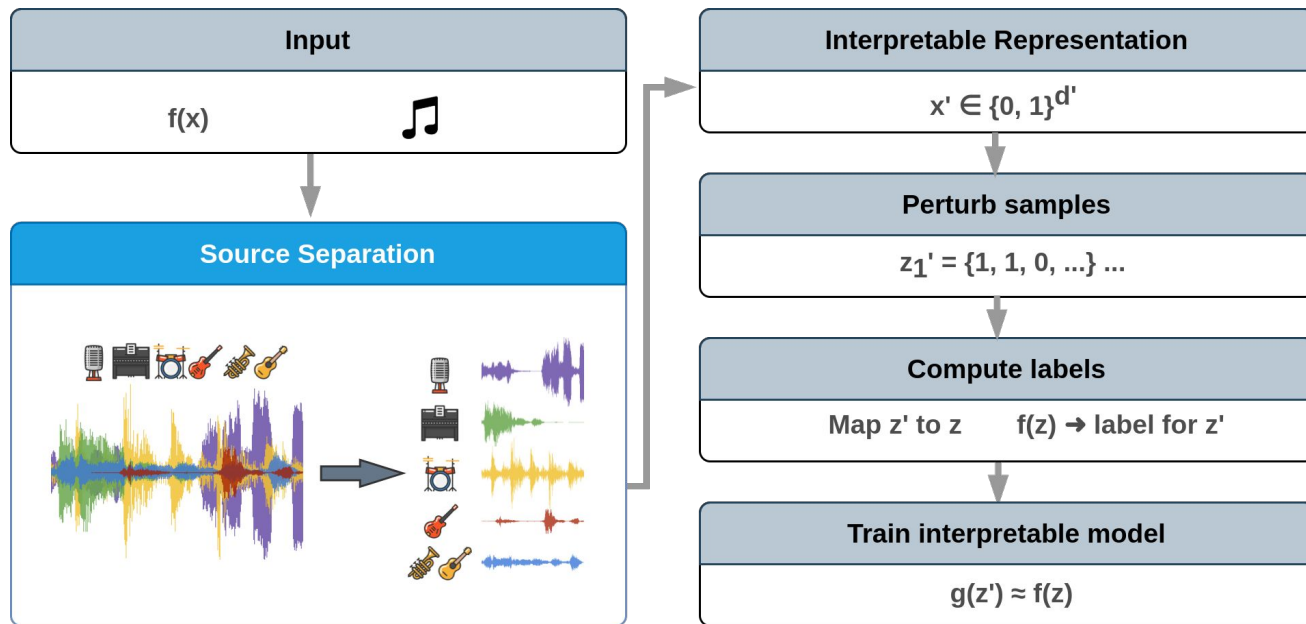
Mid-level Explanations for Valence in Hip-hop



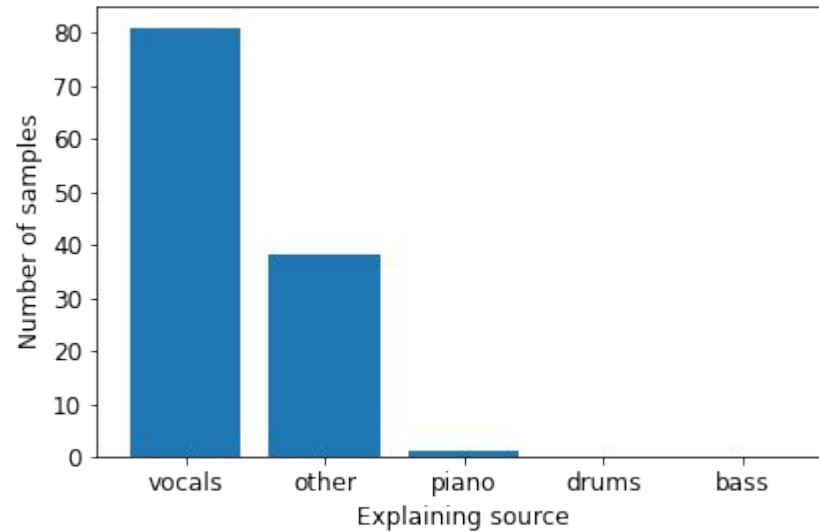
What musical component is the “cause” for high rhythm stability predictions?



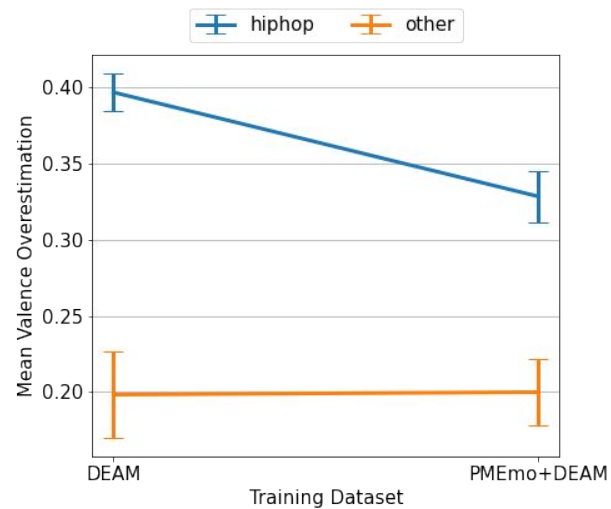
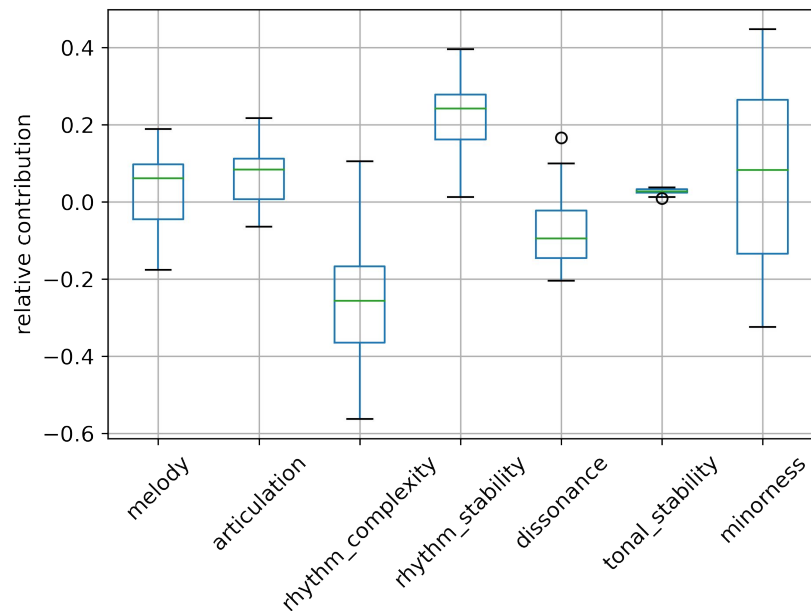
Obtaining Sound Source Explanations: audioLIME



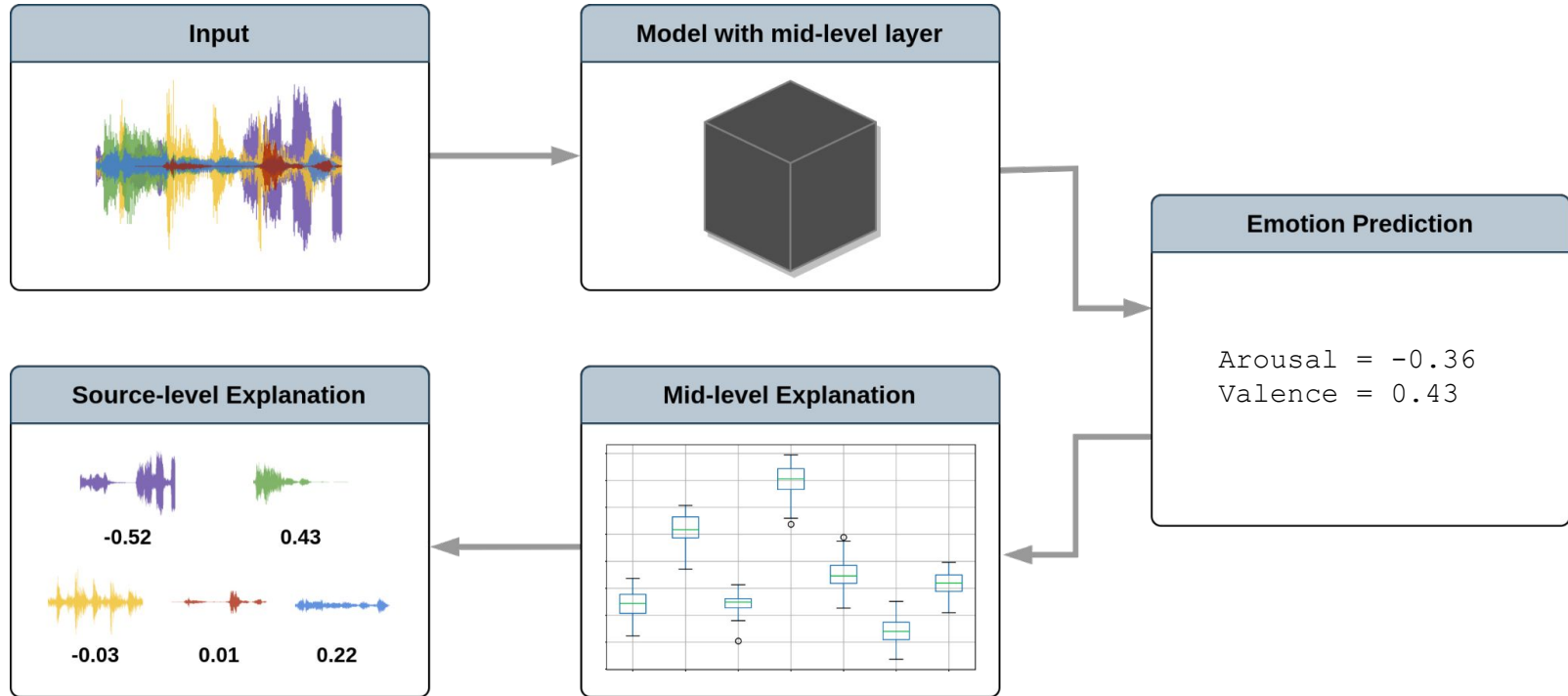
Sound Source Explanations for Rhythm Stability



Retraining with Combined Training Set



Schematic



Schematic

