# Ilaria Manco

Emmanouil Benetos, George Fazekas, Elio Quinton

# Bridging audio and language for music understanding

Joint audio-language learning leads to better music representations and enables new music-related tasks.

Are we doing joint multimodal training right?