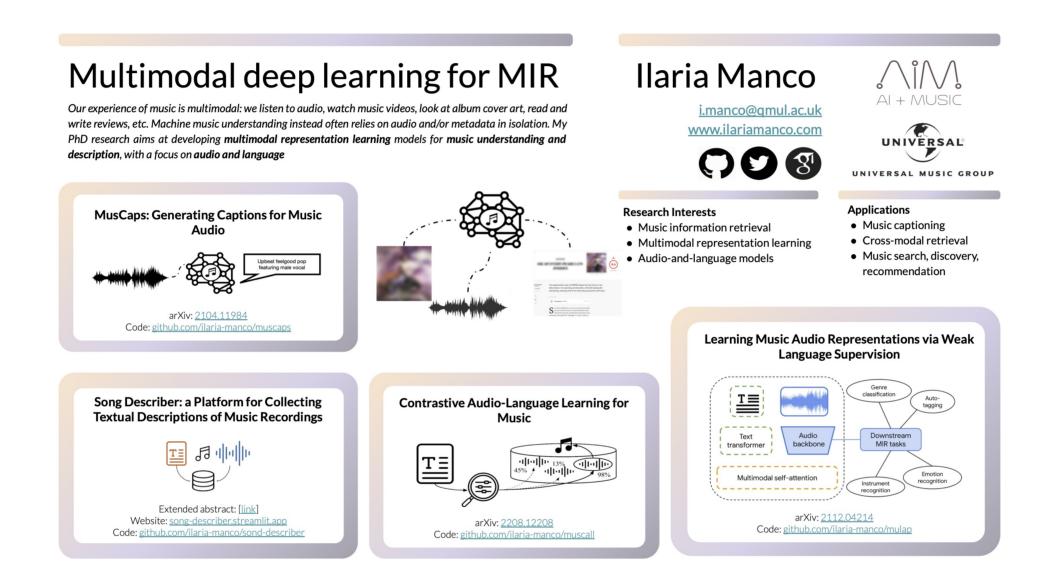# Bridging audio and language for music understanding

## Ilaria Manco

2019



# Finding

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

# Question

Are we doing joint multimodal training right?

Supervisor(s): Emmanouil Benetos, George Fazekas, Elio Quinton