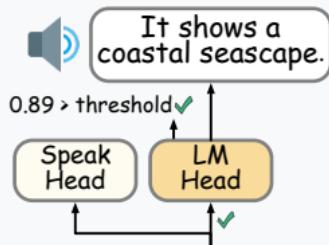
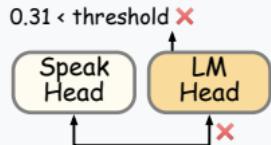
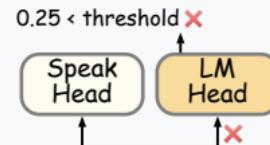
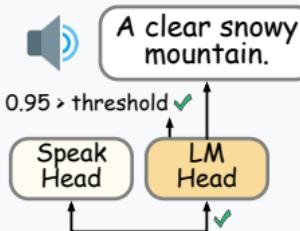


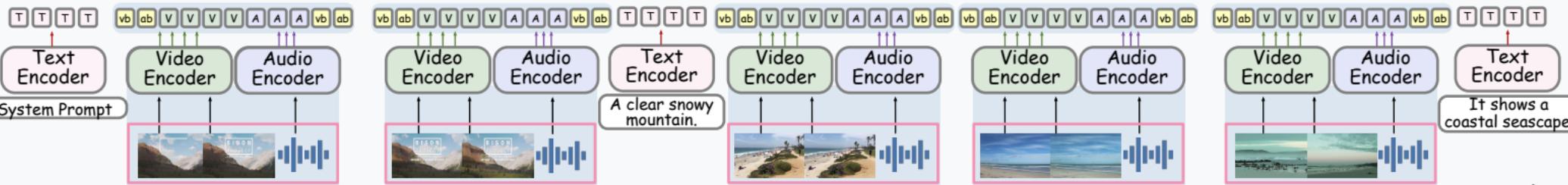


ROMA's Model

- Multimodal Unit
- Video Stream
- Audio Stream



Omni-Modality Streaming Backbone



(How is the background in real time?)



t_1

t_n

t_{n+1}

t_{n+2}

t_{n+2+m}