# Evaluation: IMSM

9/27/2024, Thanmaya, Mus2Vid

# Overview of progress on tasks

- Read the MeLFusion Paper especially their evaluation metric
- Worked through IMSM implementation basing my code off of the MeLFusion implementation

# IMSM (Image Music Similarity Metric)

- MeLFusion discusses a Multimodal approach to generating audio music from text and image
- Utilizes CLIP and CLAP scores to get a music and image similarity

$$\mathcal{A}_{\text{IMSM}} = \mathcal{A}_{\text{CLIP}} \, \mathcal{A}_{\text{CLAP}}^{T}$$

# IMSM implementation

```python
20    # CLIP embeddings (Image and Text)
21    inputs = clip_processor(text=[text], images=[image], return_tensors="pt", padding=True)
22    clip_outputs = clip_model(**inputs)
23    image_embeds = clip_outputs.image_embeds
24    text_embeds_clip = clip_outputs.text_embeds
25
26    # CLAP embeddings (Audio and Text)
27    inputs_audio = clap_processor(audios=[audio], text=[text], return_tensors="pt", padding=True)
28    clap_outputs = clap_model(**inputs_audio)
29    audio_embeds = clap_outputs.audio_embeds
30    text_embeds_clap = clap_outputs.text_embeds
31
32    # Compute cosine similarities between embeddings
33    cos_sim_clip = torch.nn.functional.cosine_similarity(image_embeds, text_embeds_clip)
34    cos_sim_clap = torch.nn.functional.cosine_similarity(audio_embeds, text_embeds_clap)
35
36    # IMSM Metric Calculation
37    imsm_score = torch.matmul(cos_sim_clip, cos_sim_clap.T)
38    print(f"IMSM Score: {imsm_score.item()}")
```

# Next Steps:

- Test IMSM on MeLBench data set to see results
- Try using on CoLab, and Gilbreth to get different results