

 Freezing  Training/Fine-tuning

Agent: The video is about a cat
and ..., with a ... music

User:
<visual-tokens>\n
<audio tokens>\n
Can you describe
the video and the
audio?

Visual Tokens

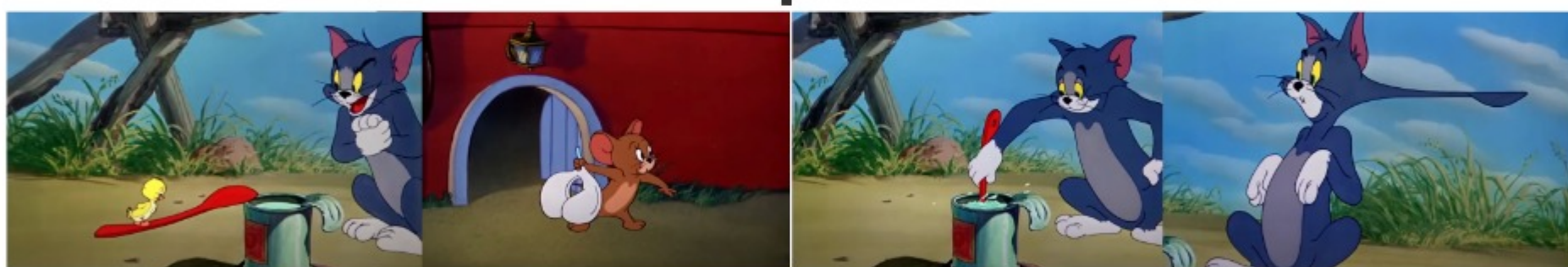


Visual Projection Layer



Spatial Token Merging

Visual Encoder
(CLIP)



Event-based sampling



Input Video



Audio Projection Layer



Audio Encoder
(MERT)



Audio Tokens



Large Language Model (Vicuna 7B)

