Frame-level Multimodal Representations Modality-specific Modality-invariant "I love life" Recognition **Fusion** frames Audio stream Video stream