

Video Frame Preprocessing

Feature Extraction

Multi-modal Feature Processing

Cross-modal Fusion

