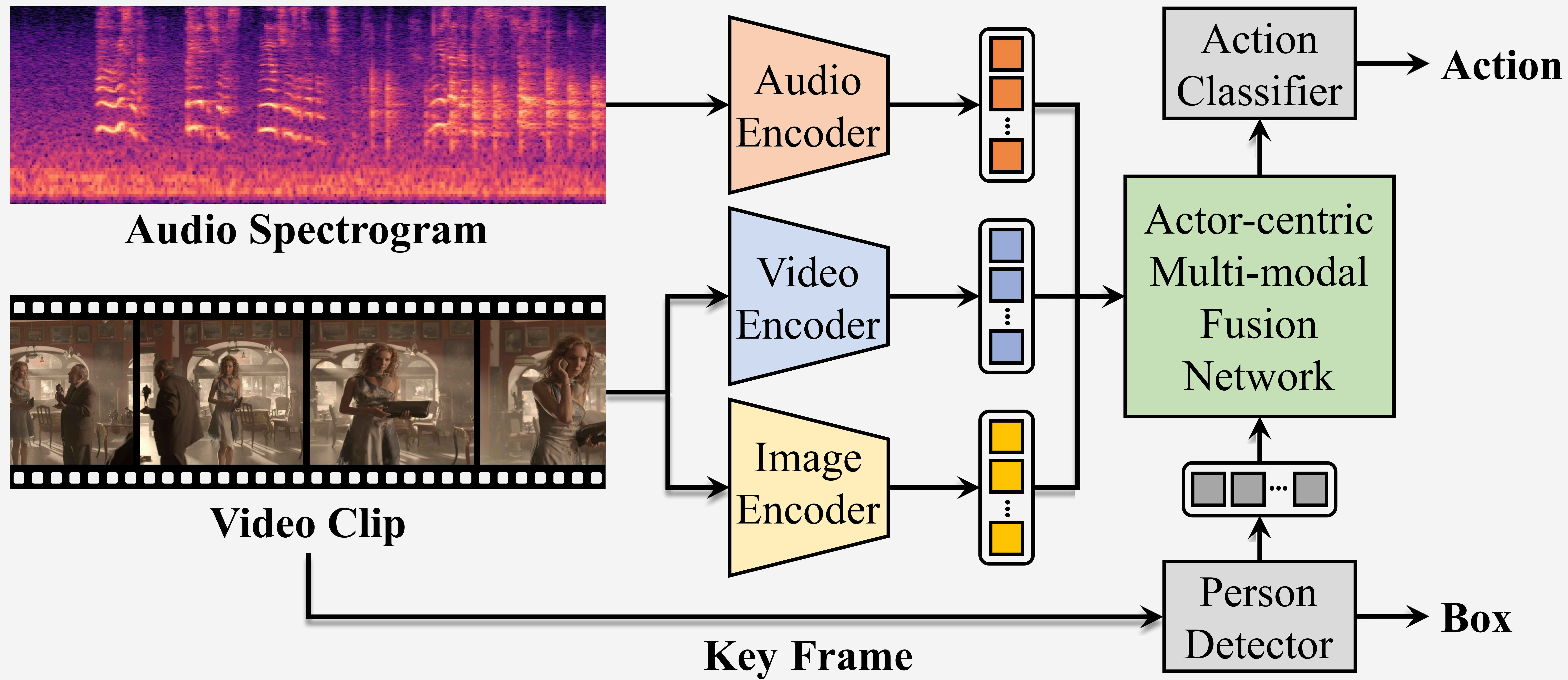
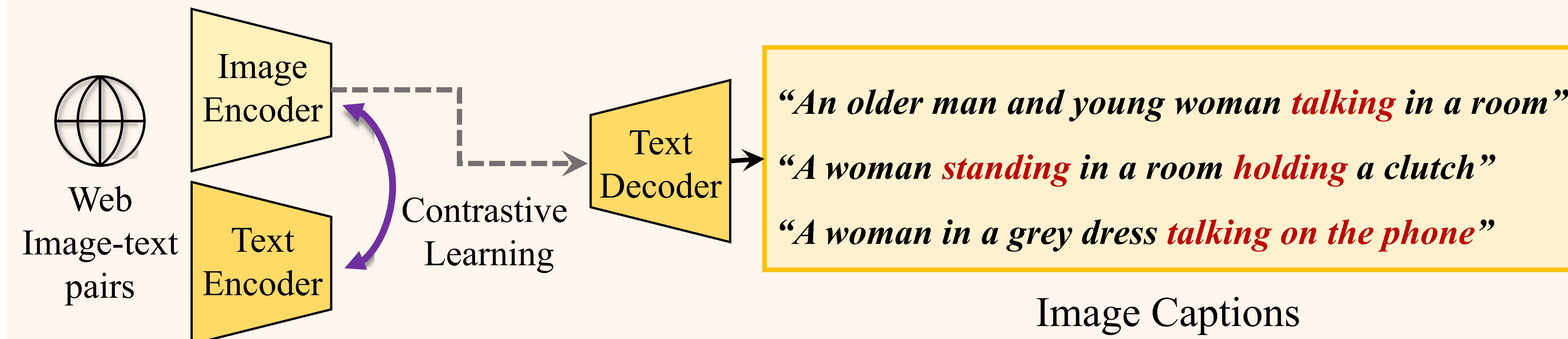


Overall Framework of JoVALE



Vision-Language Pre-training & Fine-tuning for image captioning



Actor-centric Multi-modal Fusion Network

