# Modeling Videos: Language as a Key Driver

Research Seminar: Video–Language Models[1]

Lei Wang

Australian National University

December 5, 2023

Australian
National
University

# A multi-branch cross-modal pre-training framework
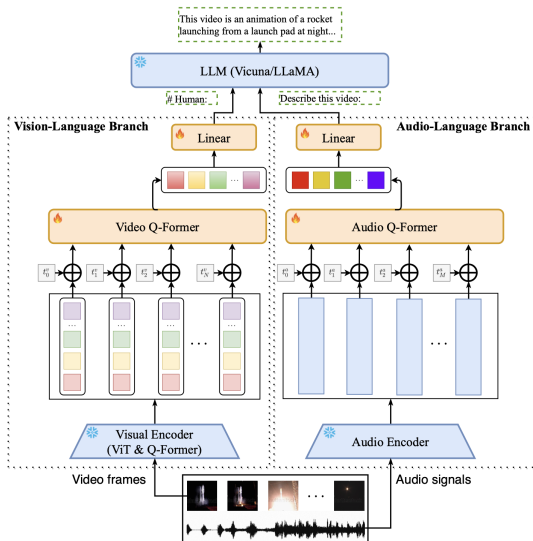


Figure 1: Video-LLaMA.

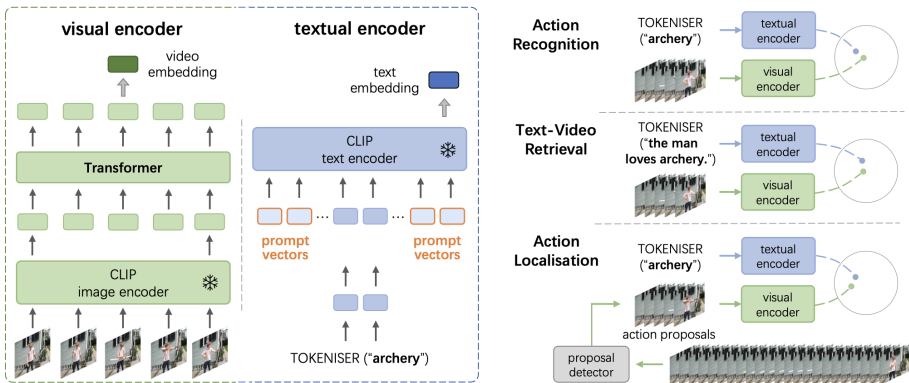# Prompting visual-language models



Figure 2: Model adaptation by learning prompts and temporal modelling.
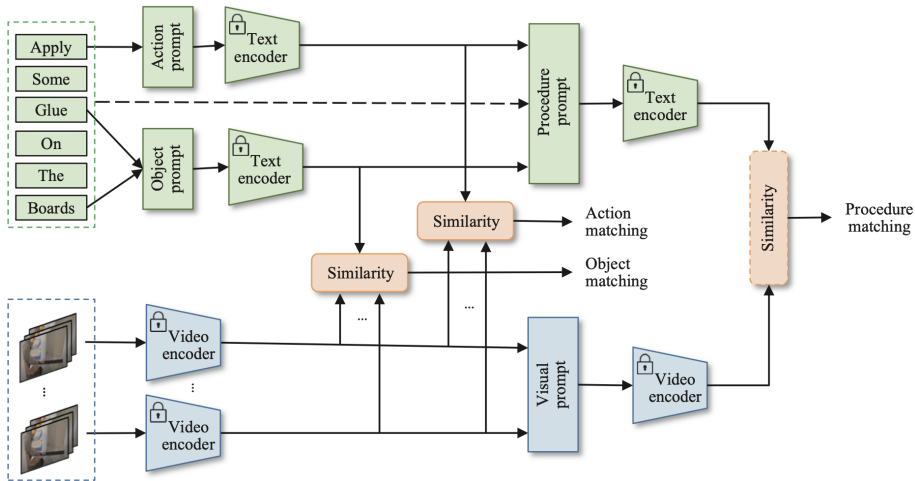
# Compositional prompt learning



Figure 3: Compositional prompting video-language model.

# Compositional prompt learning (cont.)



(a) Procedure classification

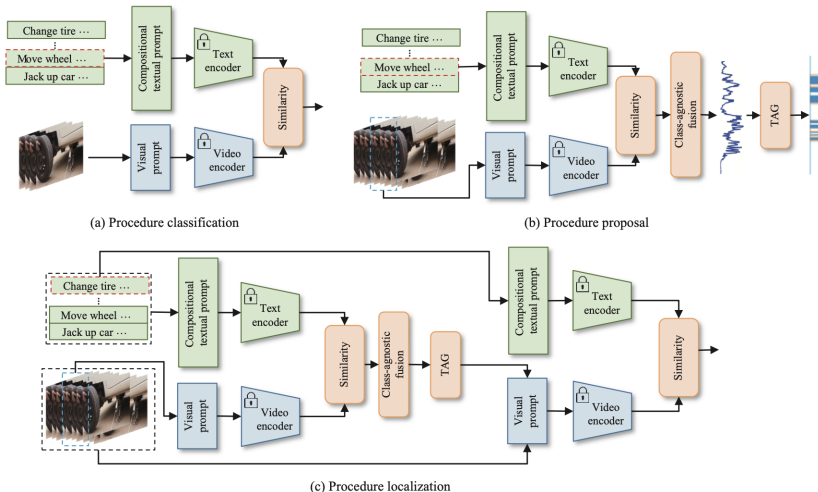(b) Procedure proposal

(c) Procedure localization

Figure 4: Three classical downstream tasks for procedure understanding are reformulated as general matching problem.

# Instilling video-language models with a sense of time



**1.**

**2.**

**3.**

**4.**

**A.** Dog runs away *before* it brings a ball to the man

**B.** The dog brings a a ball to the man *before* it runs away

**C.** The baby eats food *after* it looks into the camera

**D.** The baby looks into the camera *after* it eats food

Figure 5: Understanding the time order of events across video and language is necessary.
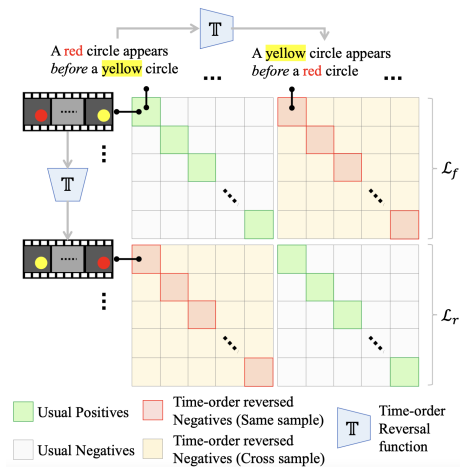
Figure 6: Negatives come from (i) other samples in the batch (cross sample) (ii) time-order reversal within the same sample.

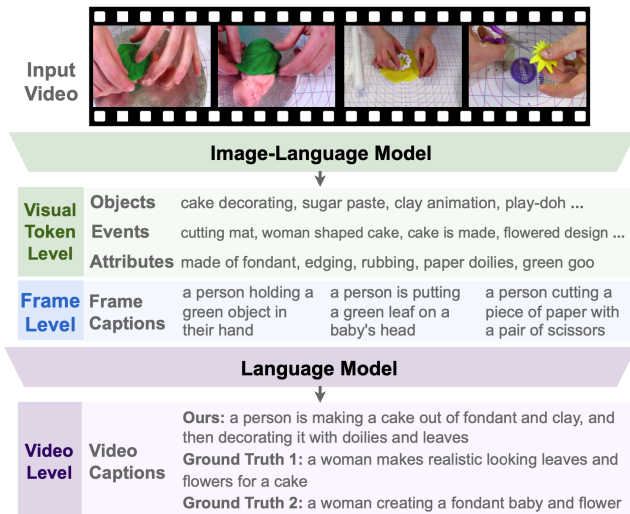# Strong video-language learners w/o pretraining/finetuning



Figure 7: Multiple levels of information in videos.
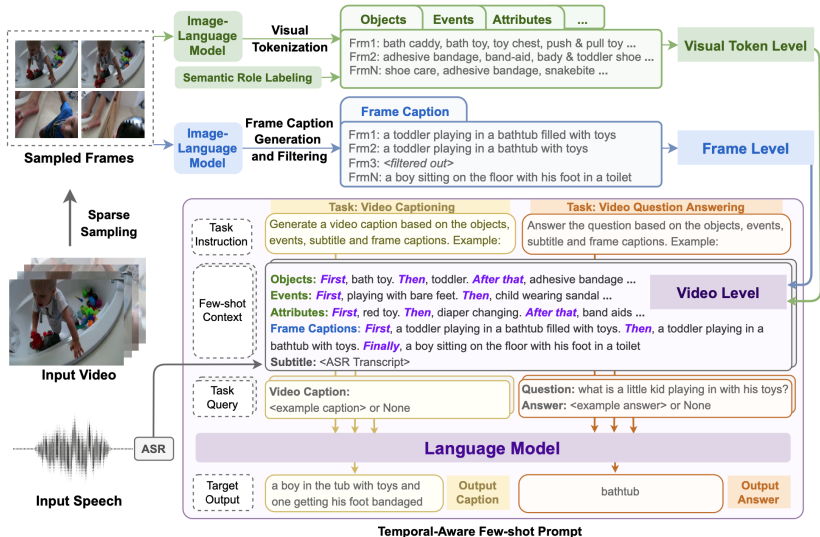
# Strong video-language learners (cont.)



Figure 8: Representing a video in a unified textural representation containing 3 semantic levels.