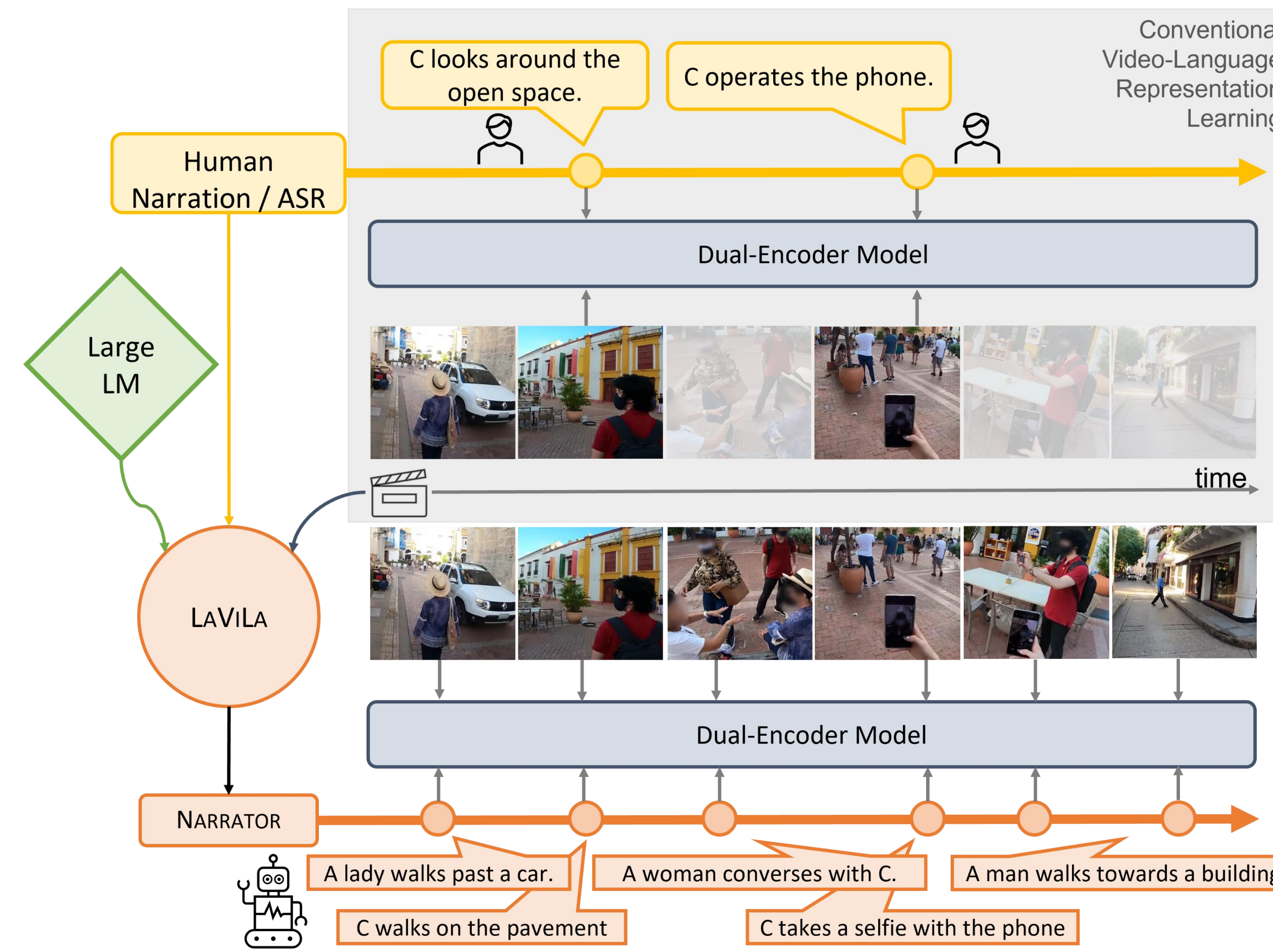
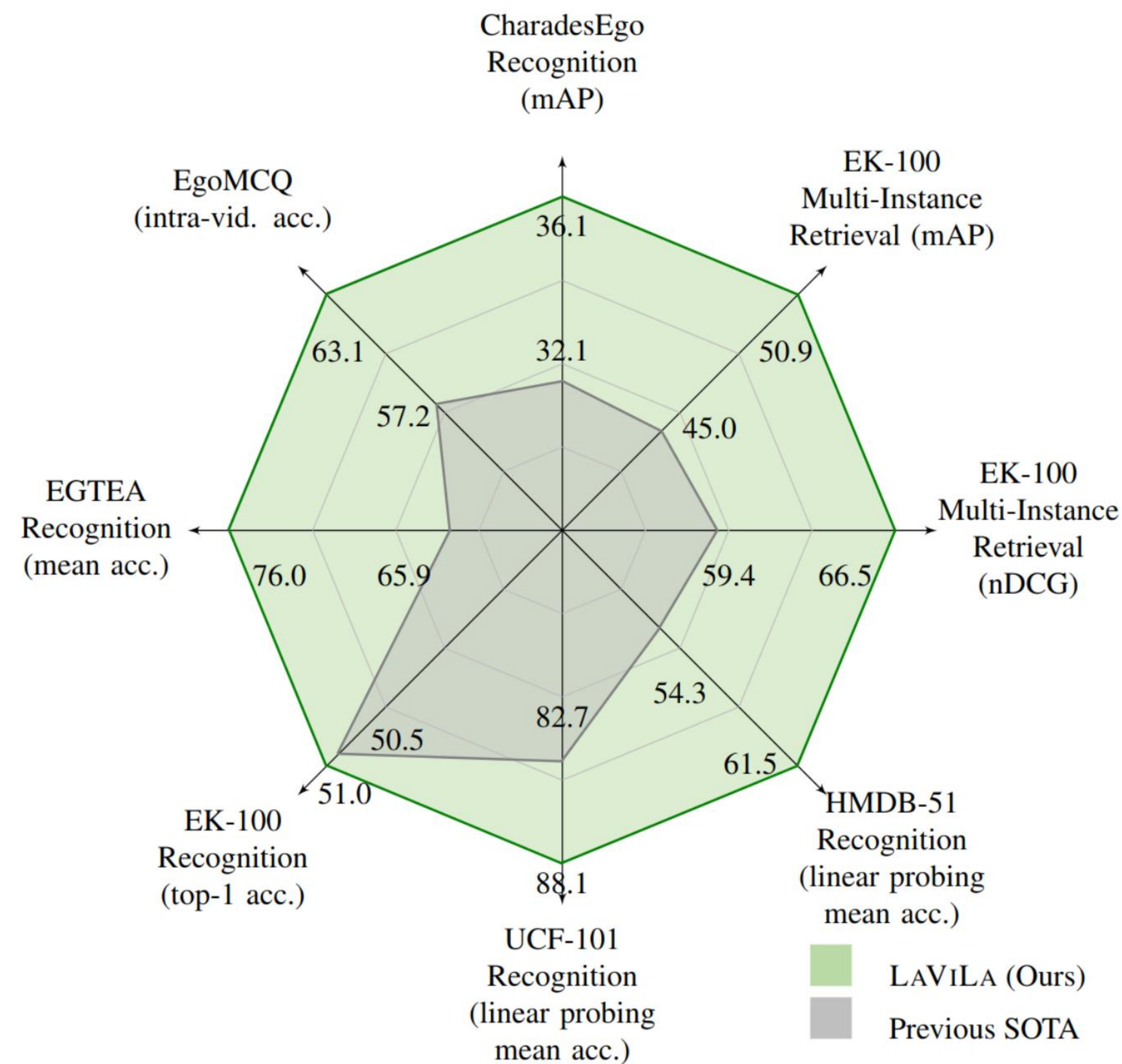
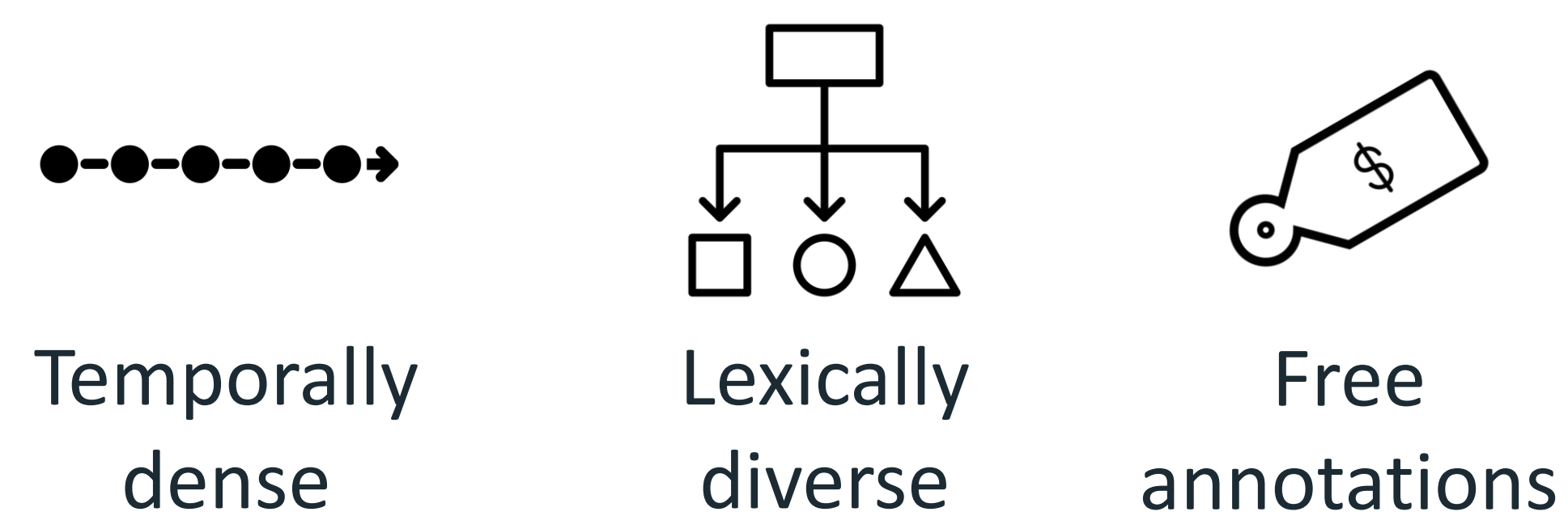




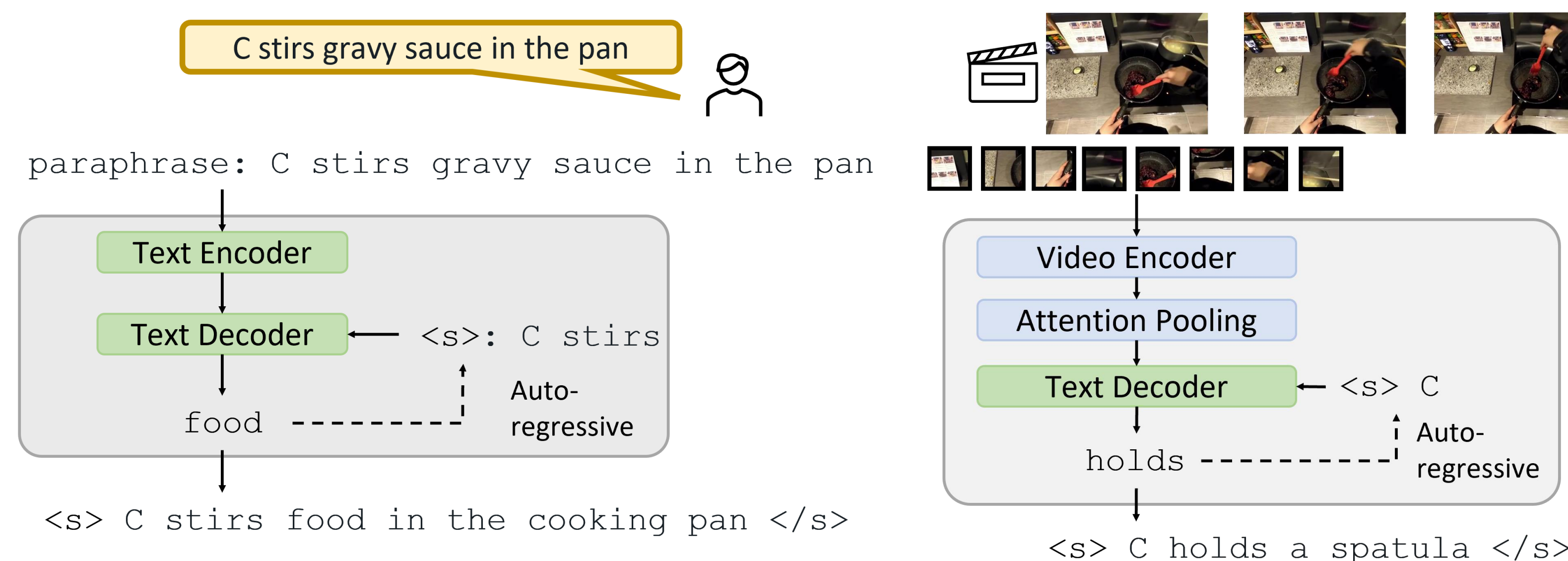
Our Approach: LAViLA

- Goal: Learn Video-Text Representation through Pre-training
- Language-model augmented Video-Language pre-training
 - LLM + visual condition => Narrators (video captioning)

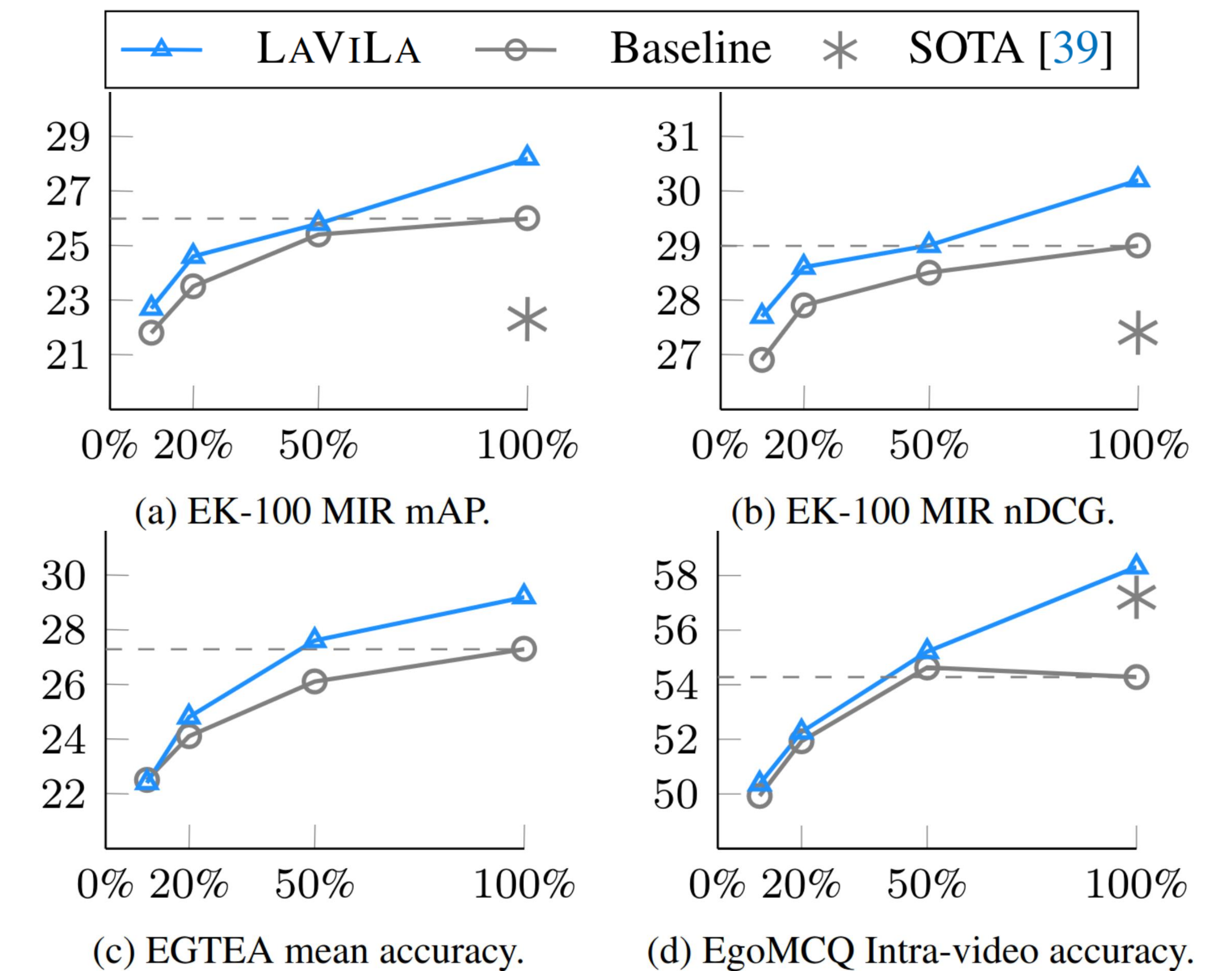


Language supervisions

- REPHRASER: narration -> rephrased narration
- NARRATOR: video clip -> pseudo-captioned narration



Data Scaling



Model Scaling

