

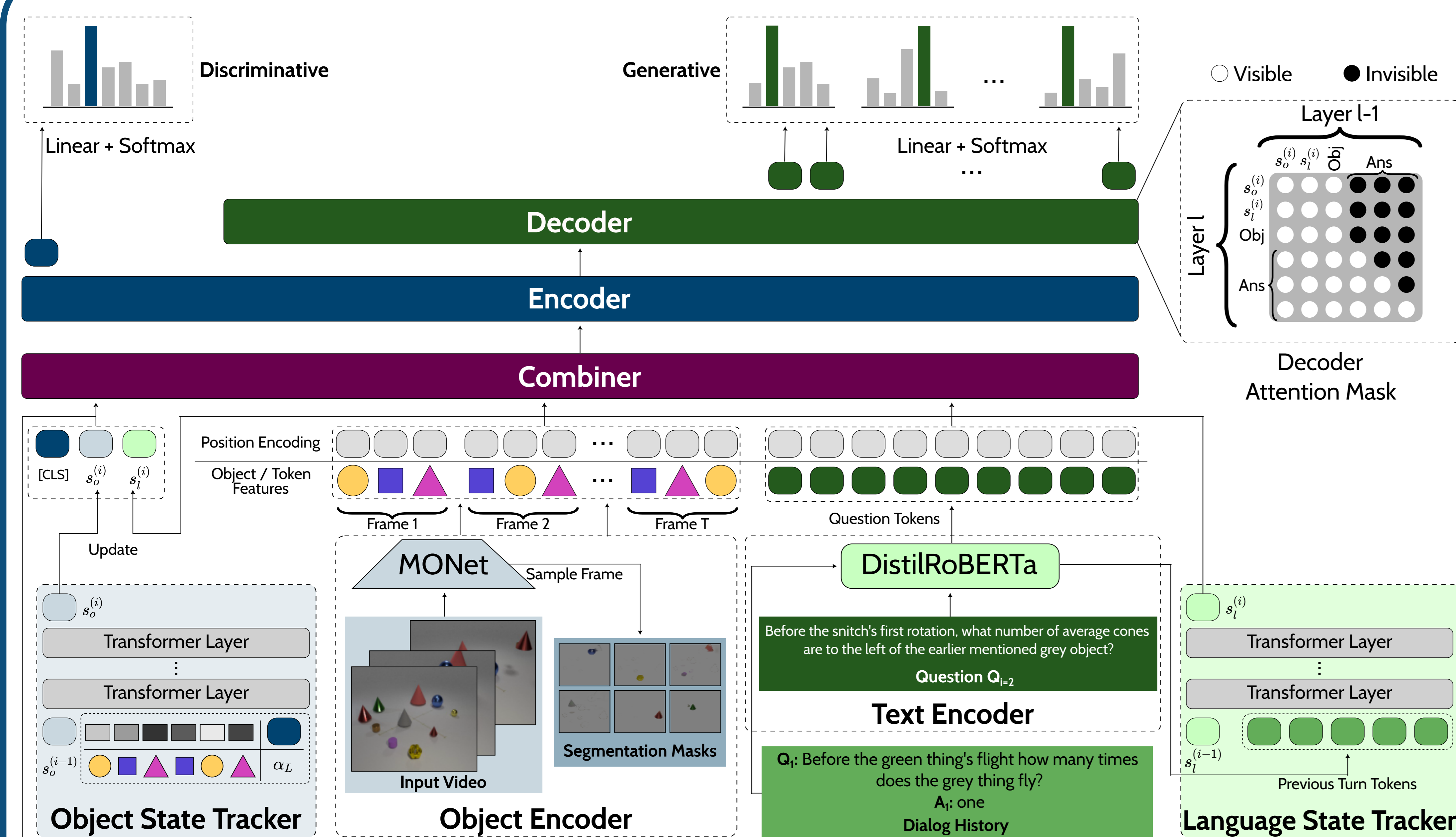
Introduction

- **Task:** Answer a question given a video and a dialog history
- **Limitations:** Prior works were trained on biased datasets that do not test higher order reasoning capabilities
- **Key insights:** Performing multi-modal state tracking alleviates the limitations of previous works and improve their reasoning capabilities

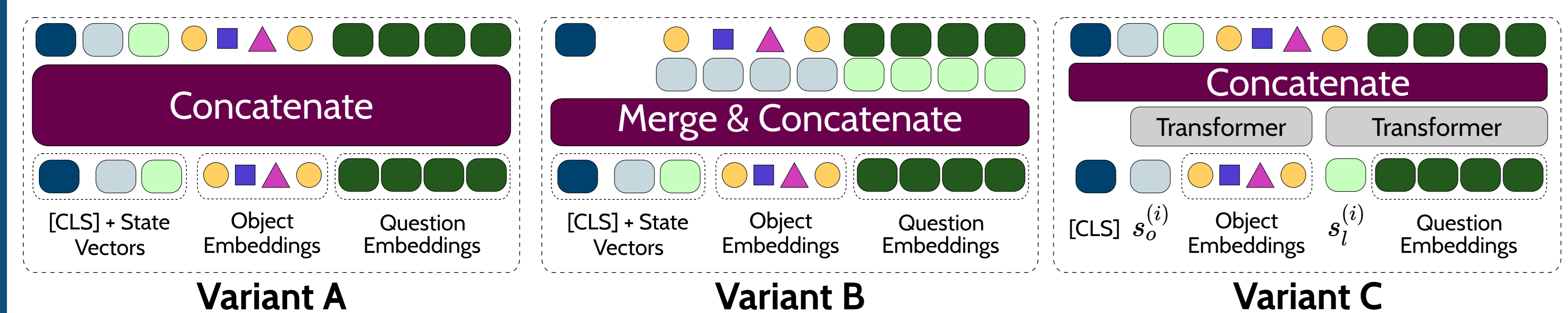
Contributions

- A novel multi-modal state tracking technique that can be seamlessly integrated into pre-trained LLMs
- Two separate attention-based trackers for object and language state tracking
- New state-of-the-art results on DVD and SIMMC datasets

Method

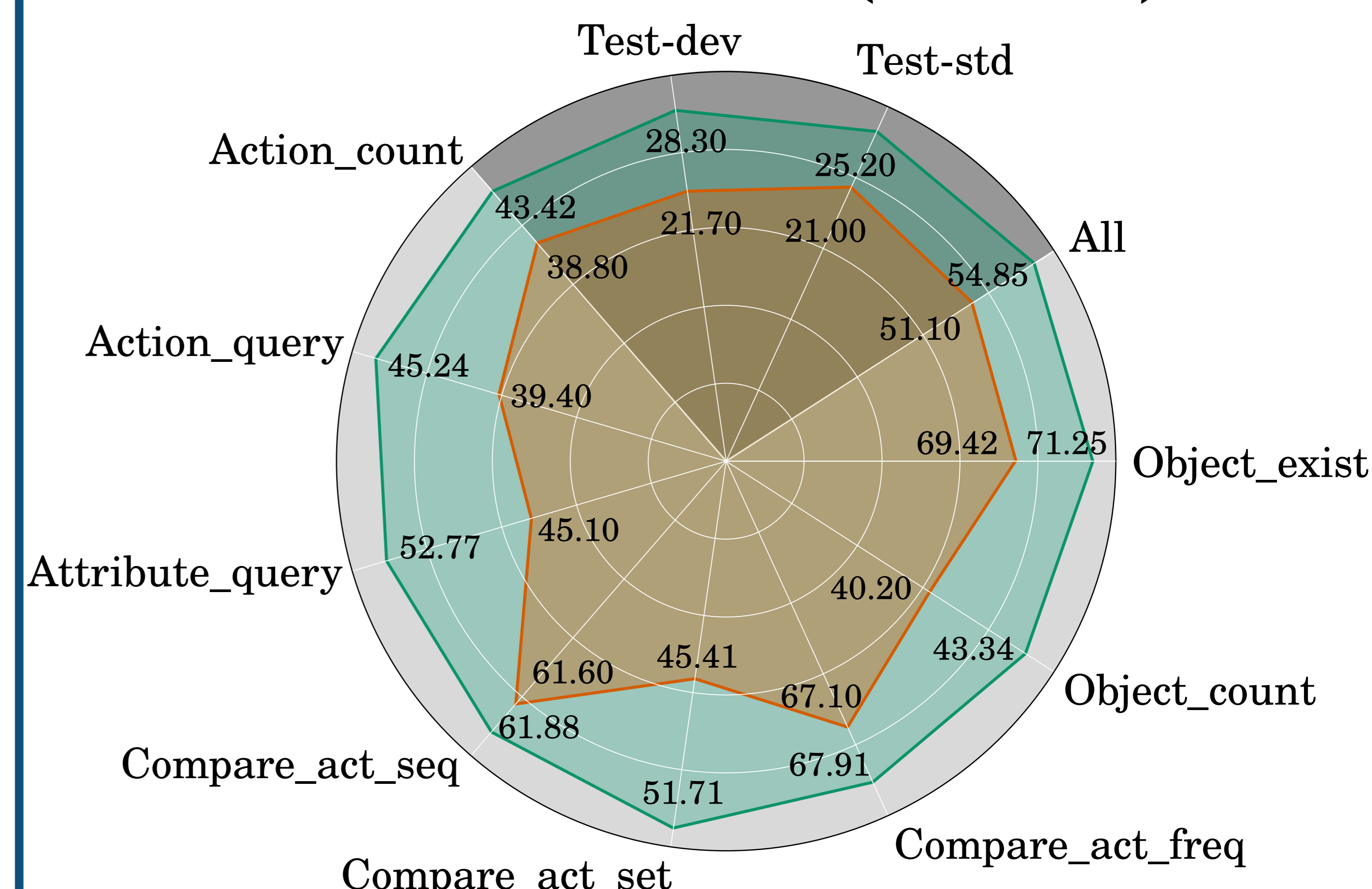


Combiner

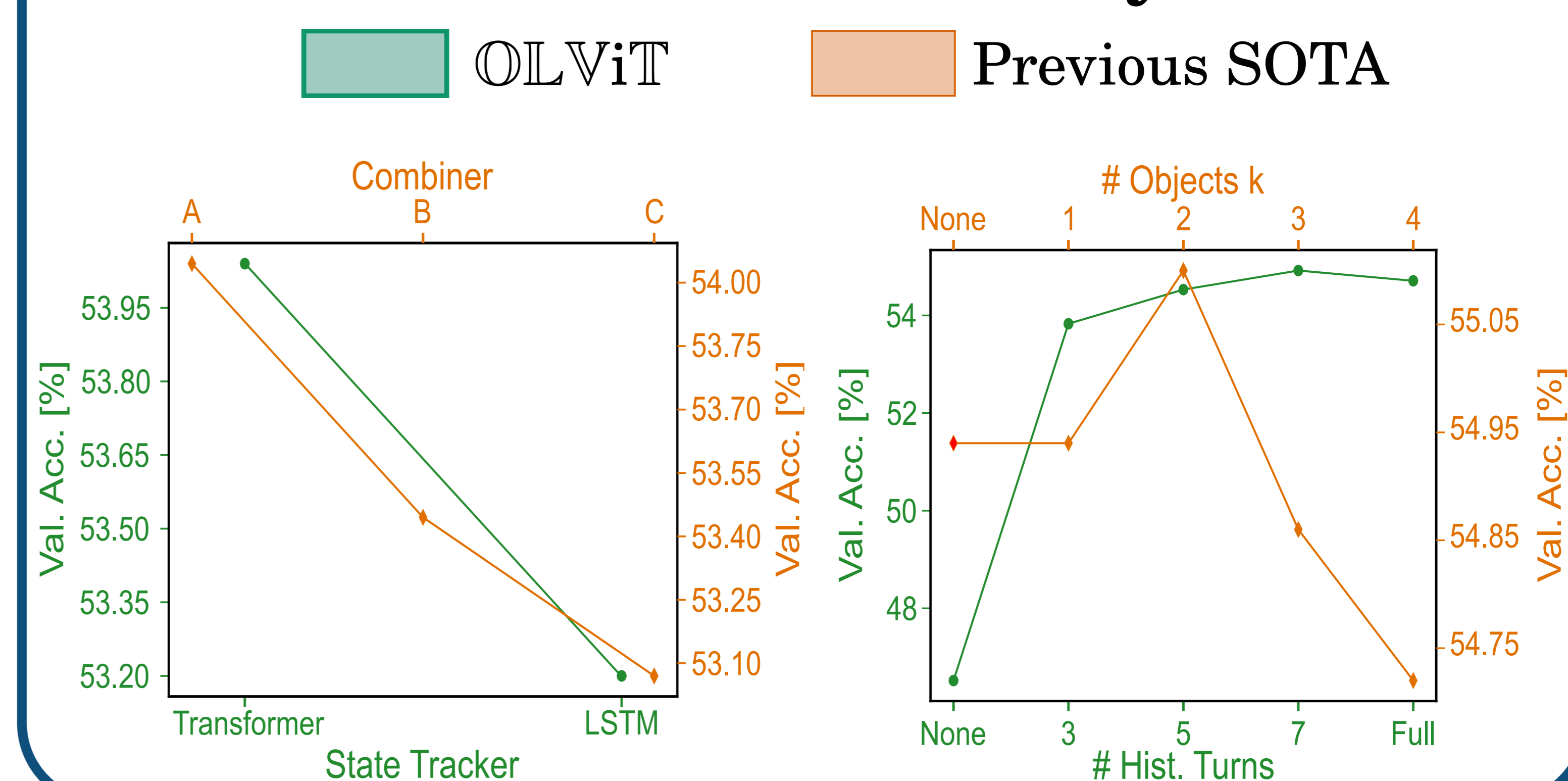


Results

SIMMC 2.1 (BLEU-4)



DVD (Accuracy)



Acknowledgments

This work received funding from the European Research Council (ERC; grant agreement 801708)