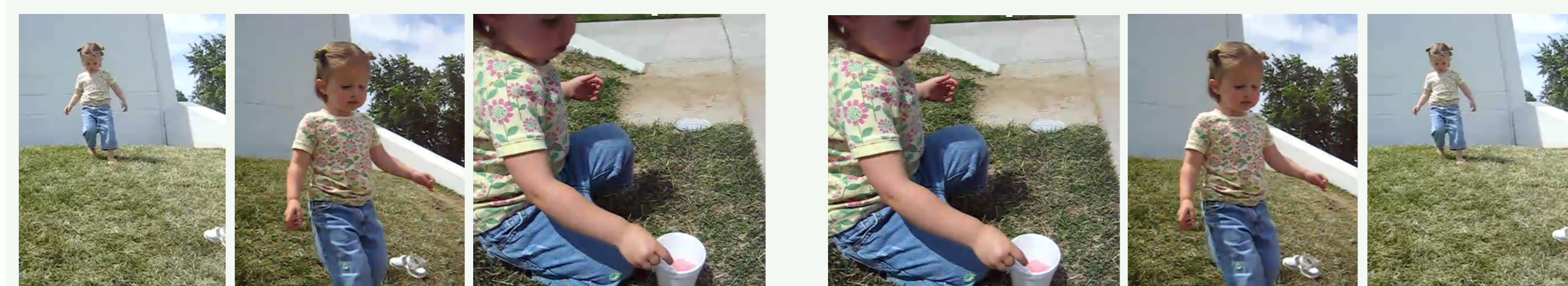




I. Can you match these video-text pairs?



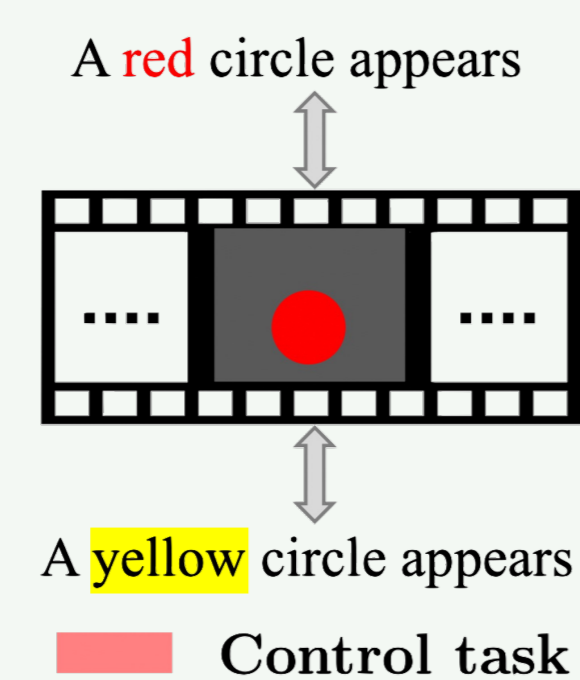
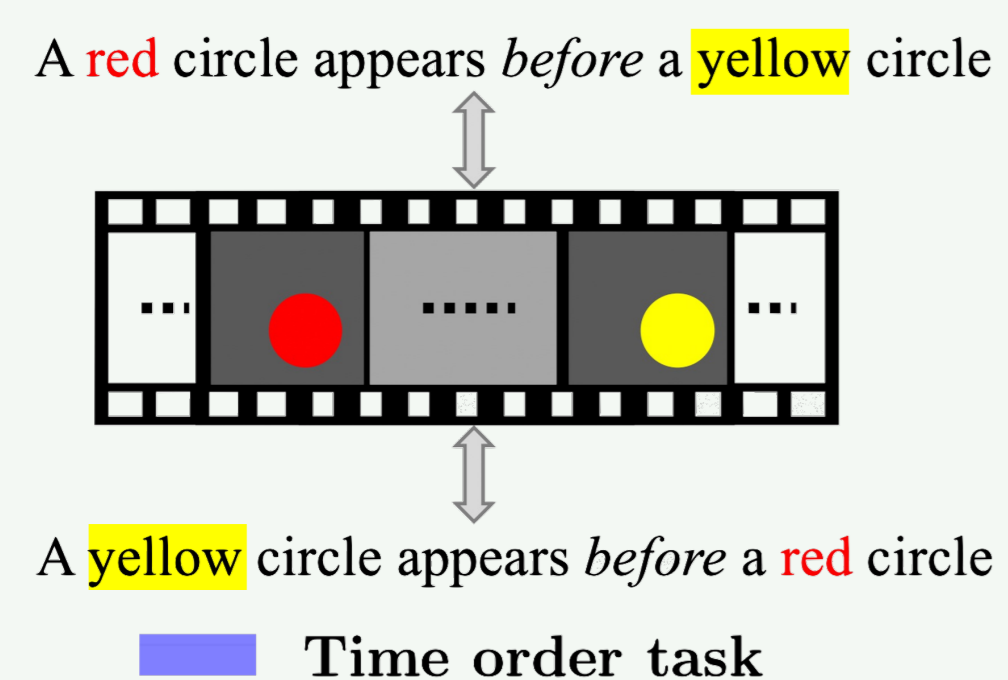
The child eats ice-cream before walking down the slope.

The child walks down the slope before eating ice-cream.

- Time-insensitive video-language models cannot match these pairs
- We propose this task as a **test of time**: matching video with its time-order consistent text

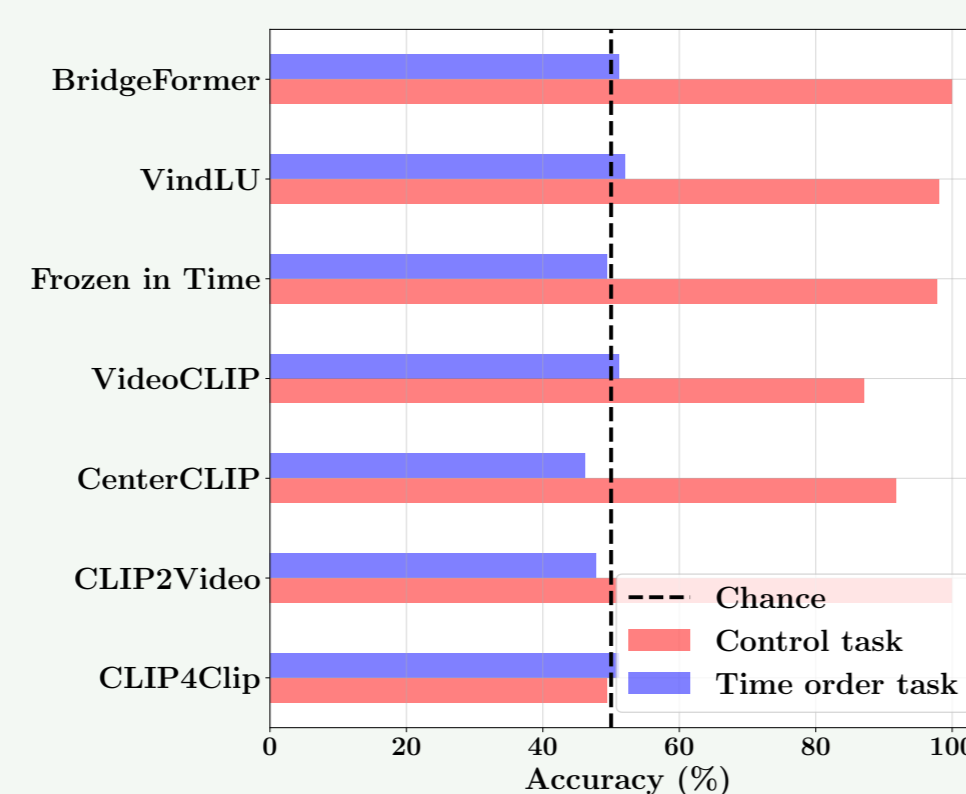
II. Do video-language models pass the test of time?

(a) Test



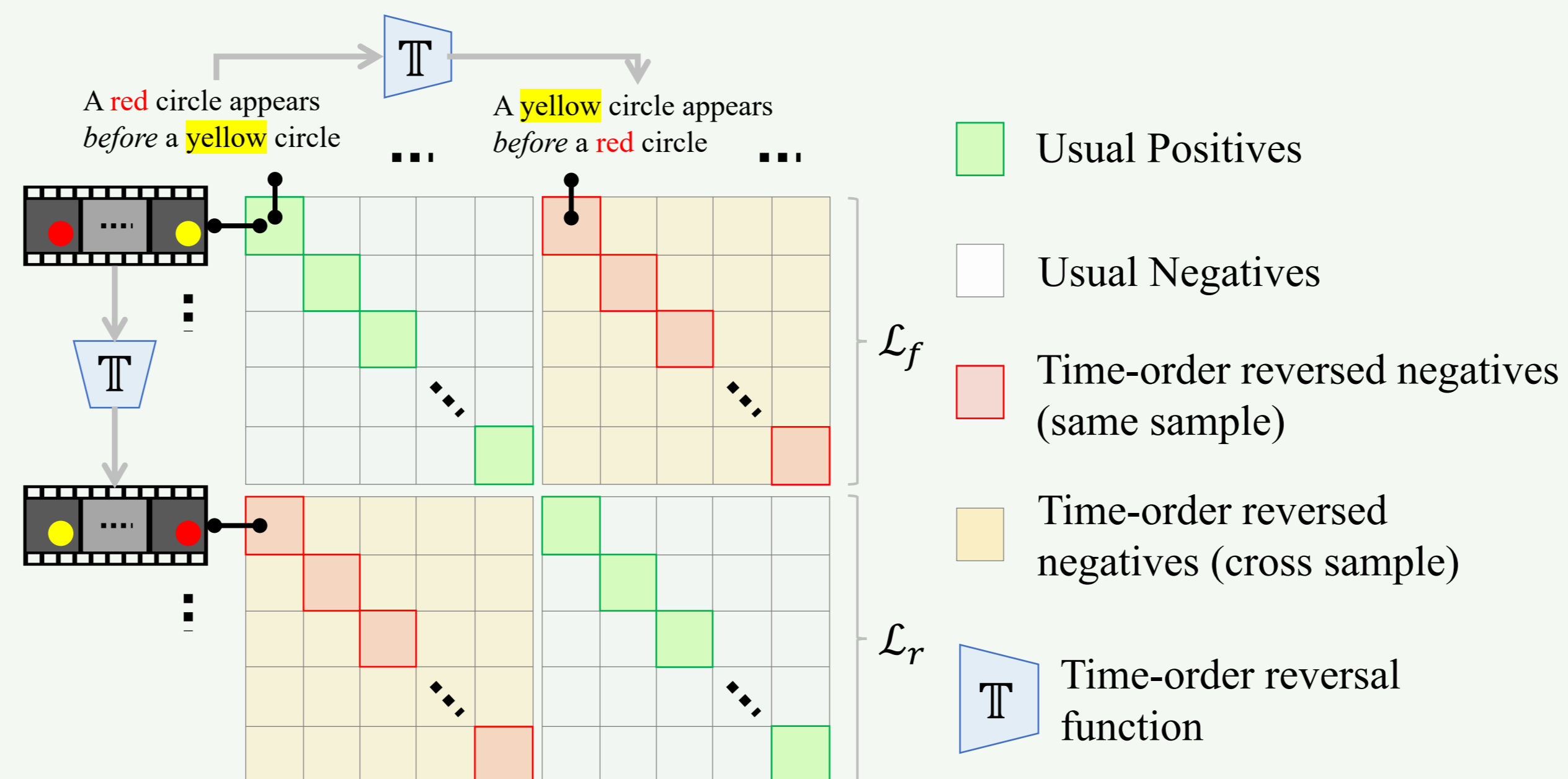
(b) Results

- We design test of time using a simple synthetic dataset
- We find that seven of the existing models fail on this test
- All excel at the control task ruling out possibility of out-of-distribution

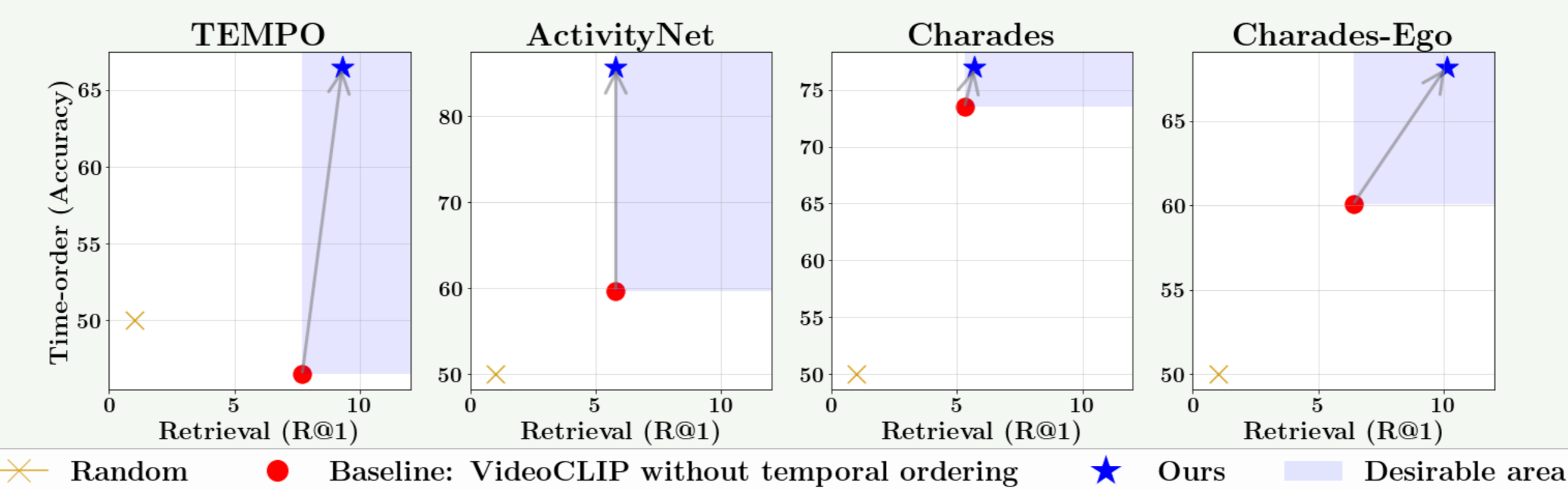


III. Can we instill this sense of time without expensive training from scratch?

- Training video-language models from scratch is super-expensive
- Our secret: carefully crafted negative samples for contrastive learning to make models time-aware
- We adapt VideoCLIP (Xu et al.) using the following recipe:

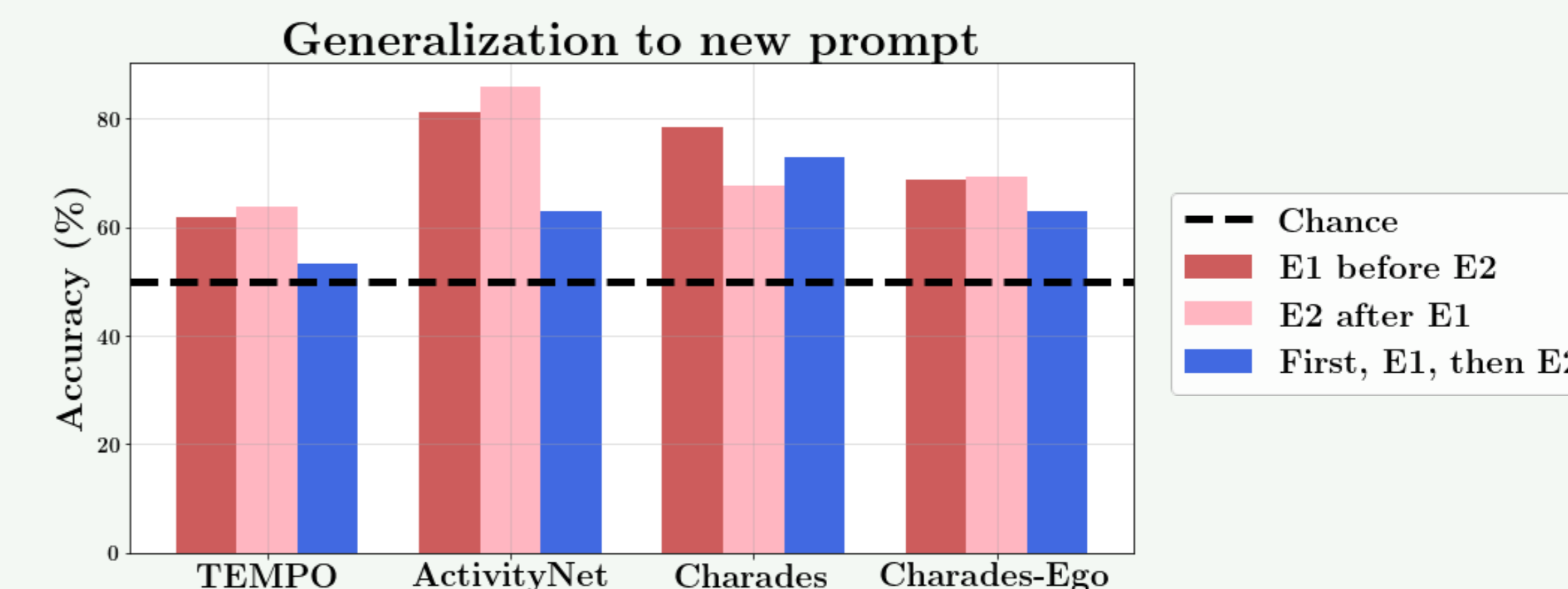


- On four diverse datasets, we show better time-awareness while preserving retrieval capabilities
- Adapted models pass test of time with flying colors! (see paper for numbers)



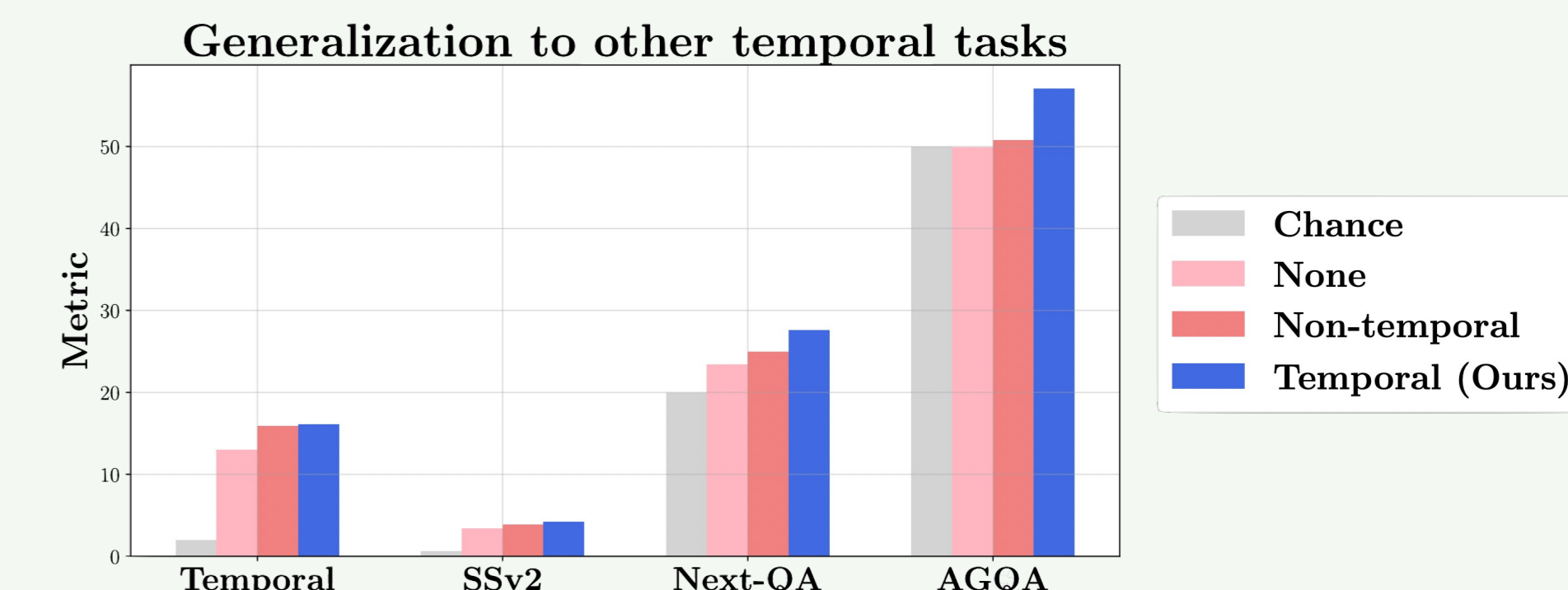
IV. Does it work beyond before/after relations?

- We test on sentences like: "First, [event 1], then, [event 2]"
- Adapted model again passes the test even with an unseen new prompt



V. Towards more general time-awareness

- Does having this narrow sense of time benefit other kinds of temporal reasoning tasks? On several temporal reasoning tasks, we observe benefits over non-temporal baseline



- We show adaptation for instilling sense of time without training from scratch. What other properties can be instilled in foundation models?