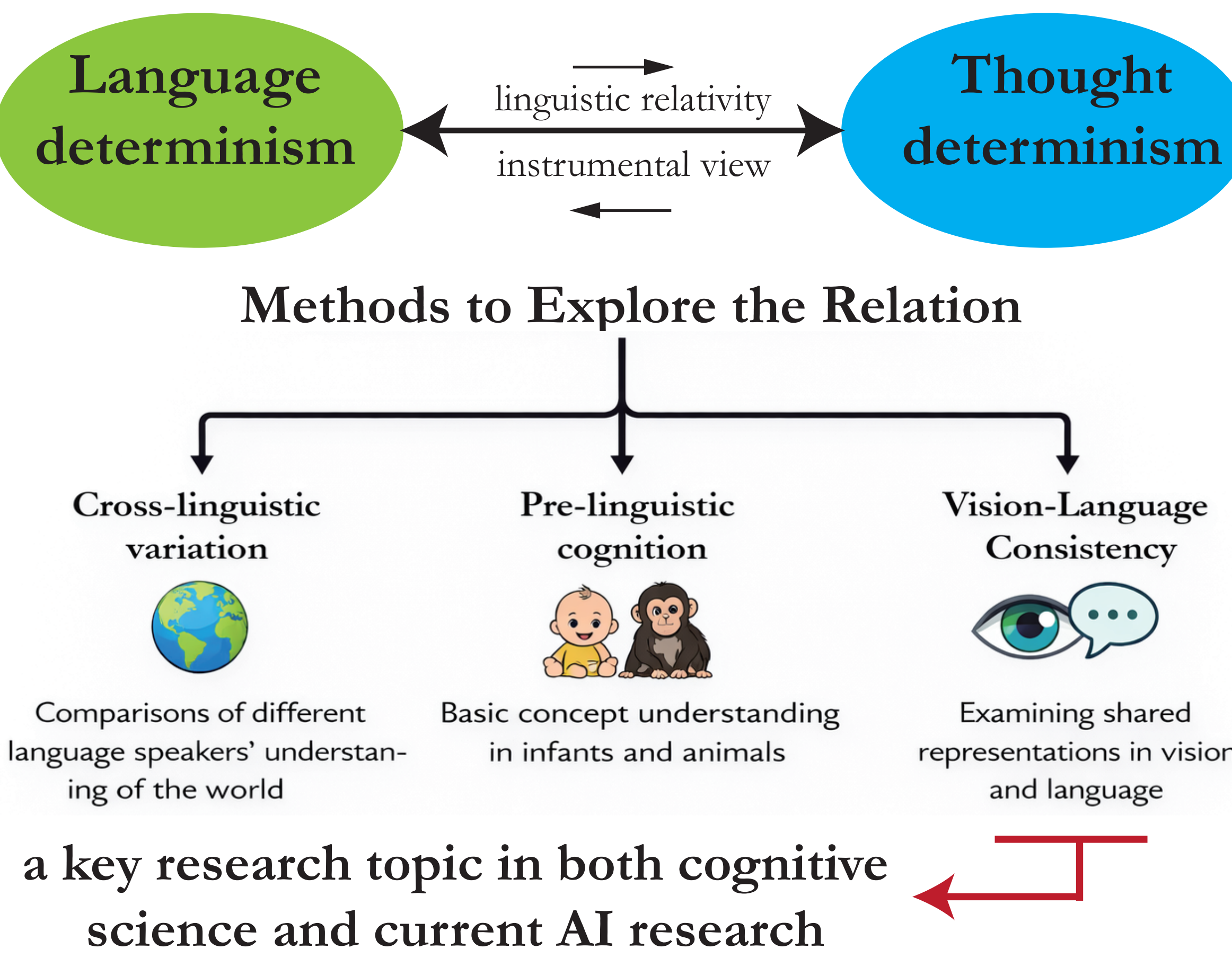
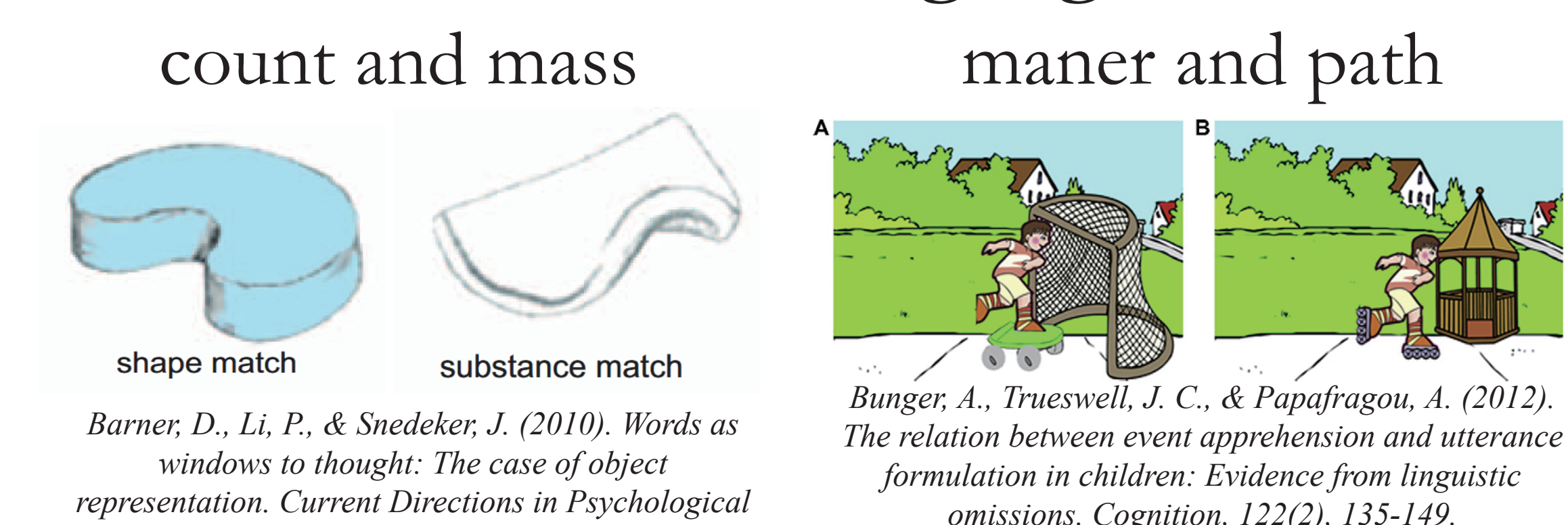


## A Longstanding Debate in Language Study

### Relation Between Language and Thought



### Parallel Structure in Language and Vision



## Grounding Verbs in Vision

### Verb Argument Structure Is Shaped by Subtle, Nonlinguistic Physical Properties

=> Certain verbs resist argument structures shared by similar verbs.

#### Example 1: Perceptual Constraint

- Someone poured the water into the glass.
- Someone filled the glass with water.
- Someone poured the glass with water. (unsuitable)
- Someone filled the water into the glass. (unsuitable)

"Pour" → directional motion  
"fill" → result state

#### Example 2: Dative Alternation

- Danielle brought the cat to her mother. ("causing-to-go")
  - Danielle brought her mother the cat. ("causing-to-have")
  - Danielle lifted the crate to him. ("causing-to-go")
  - Danielle lifted him the crate. (unsuitable)
- "lift" → continuous force do not encode transfer and caused possession.

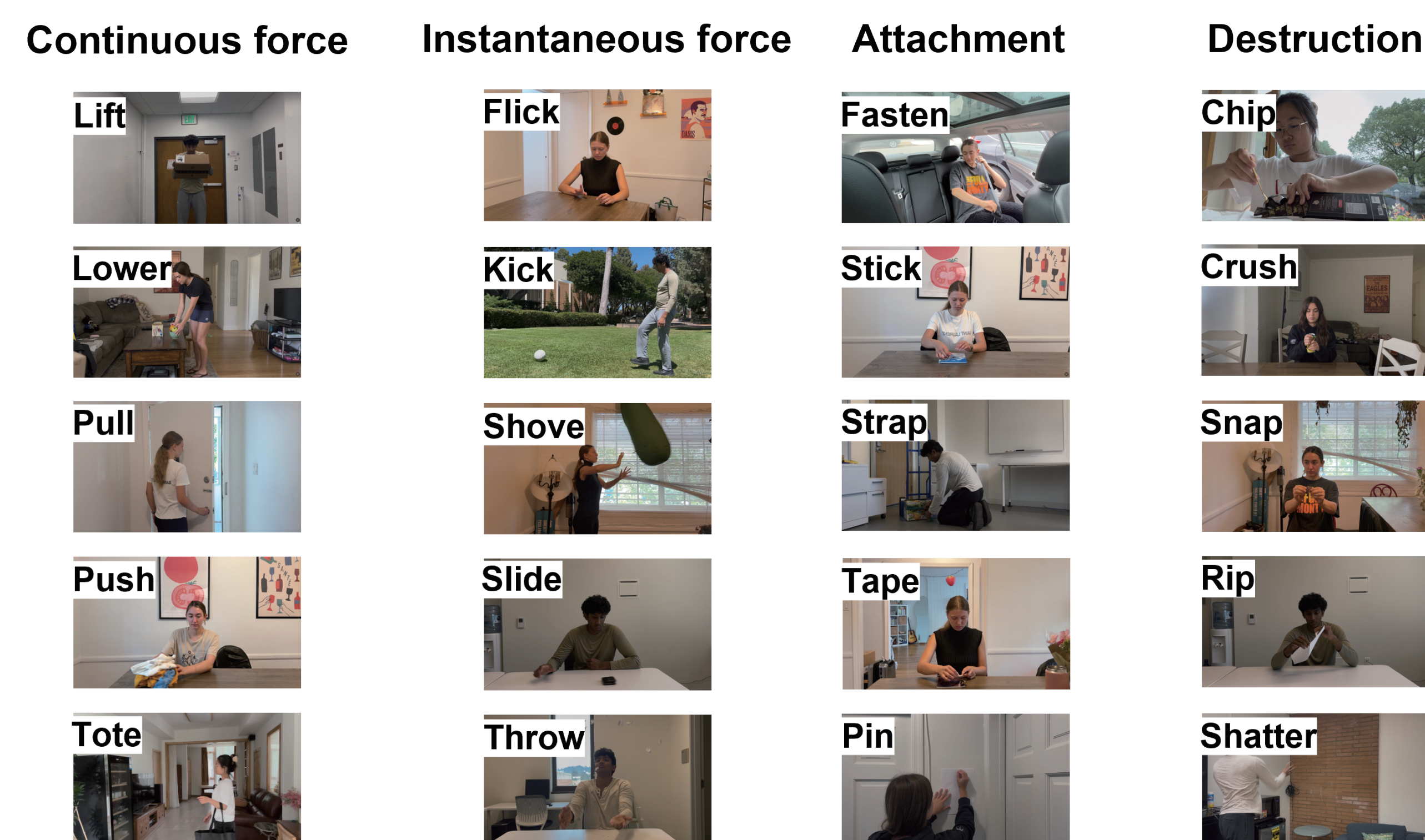
=> Humans tend to categorize these verbs based on the nonlinguistic physical causal properties.

#### Our Prediction:

Do physical causal properties that shape verb argument structure also guide few-shot learning of event categories in non-linguistic visual tasks?

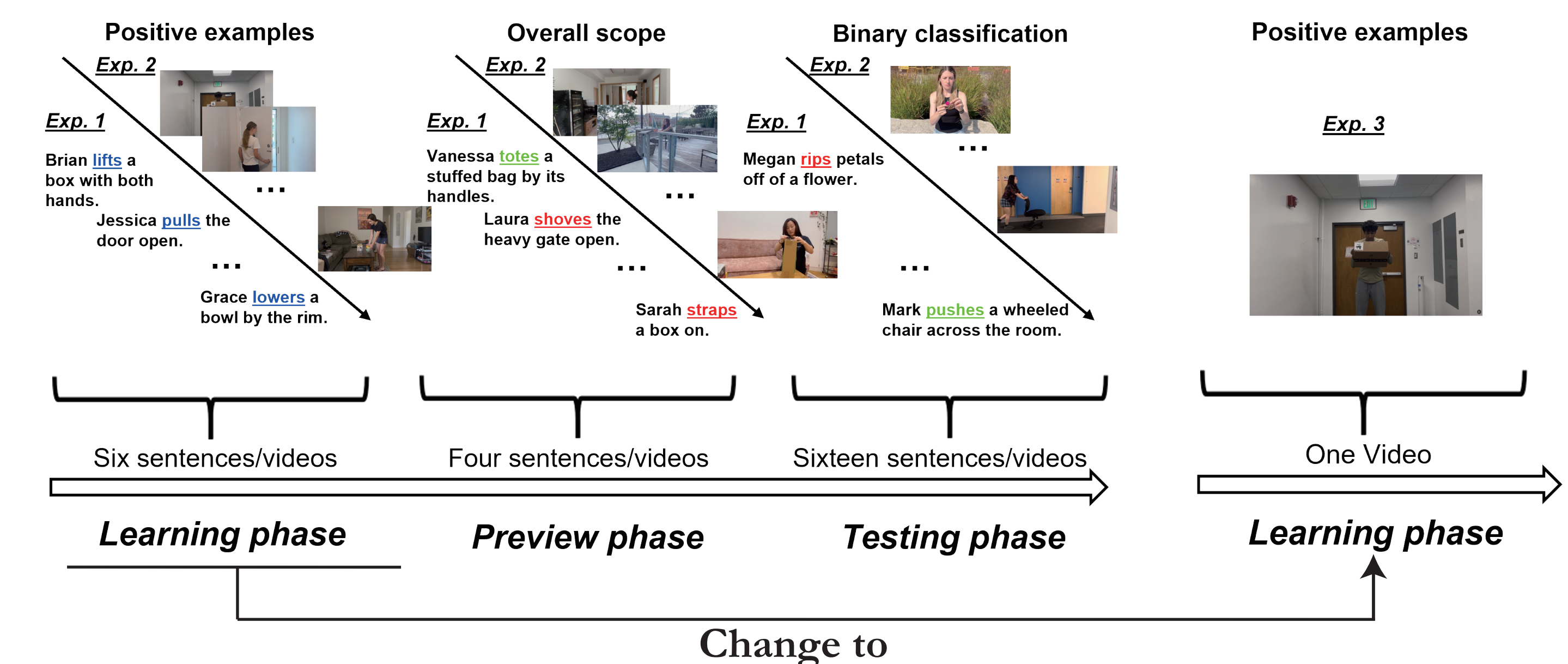
## Expt. 1-3: Few-Shot (Vision & Language) and One-Shot (Vision) Learning

### Four Classes of Verbs

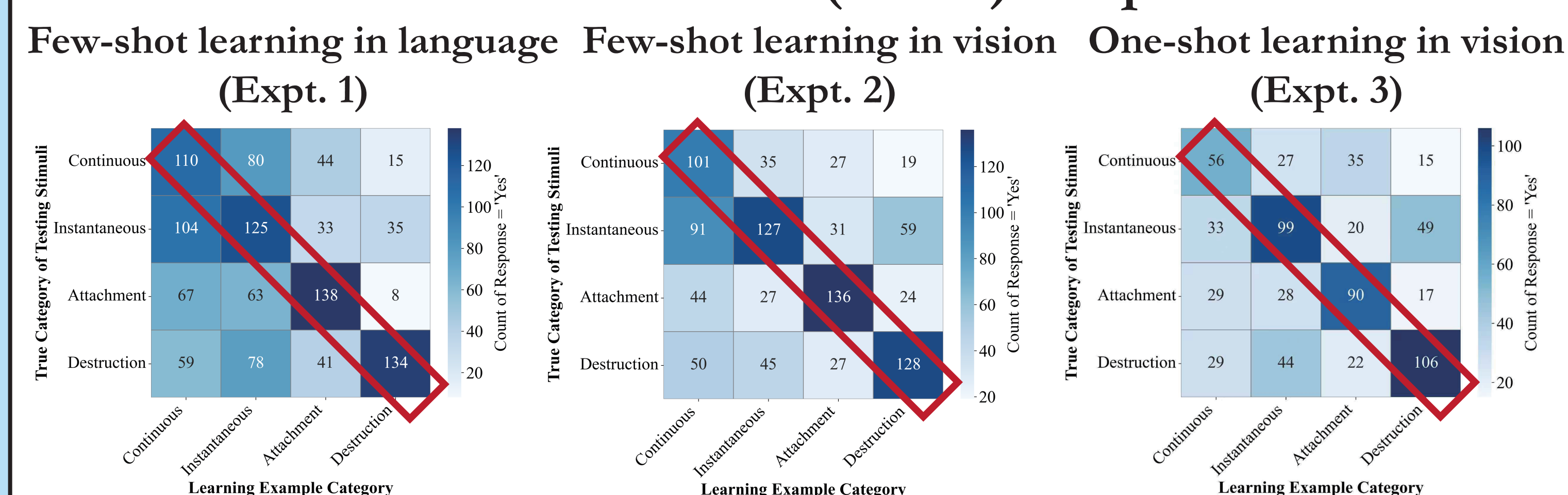


### Few-shot Learning

### One-shot Learning

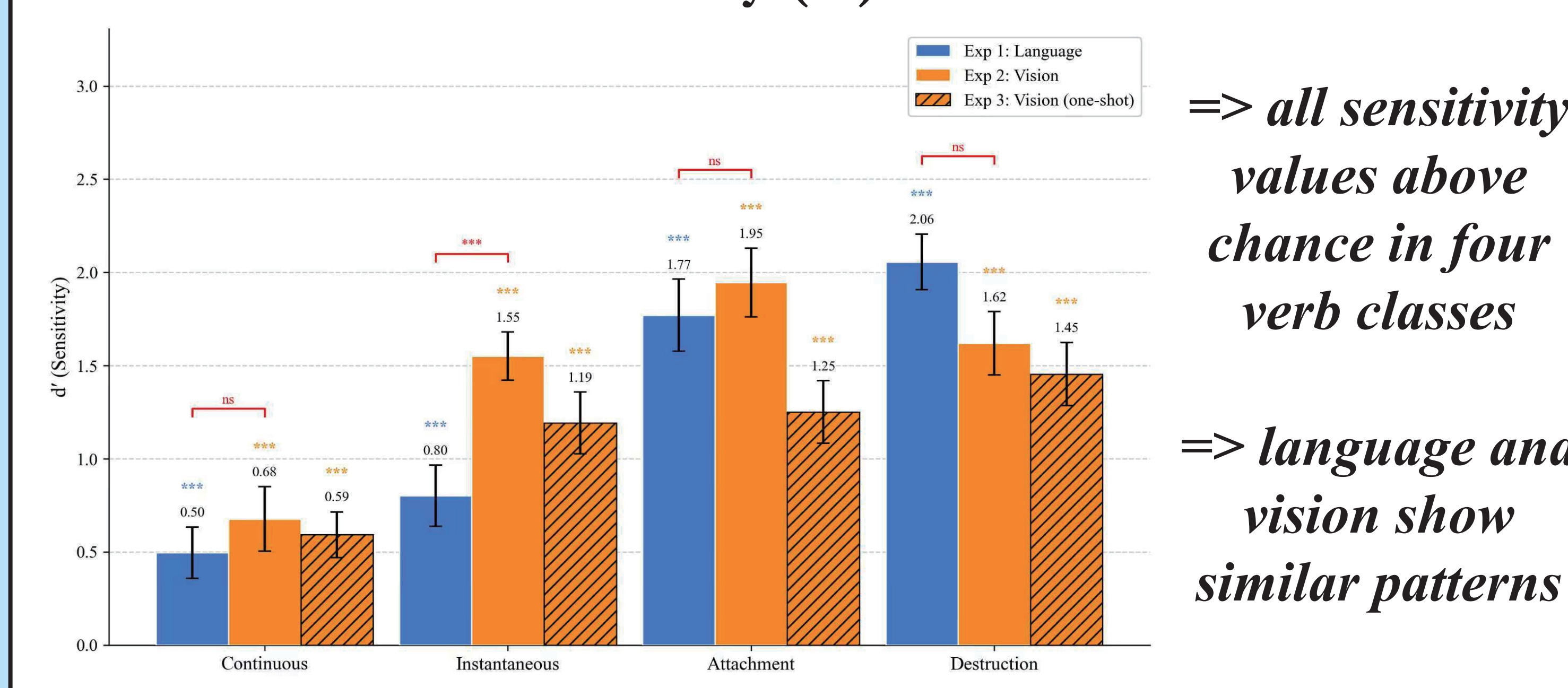


### Confusion Matrix of ("Yes") Responses



=> highest accuracy along the diagonal

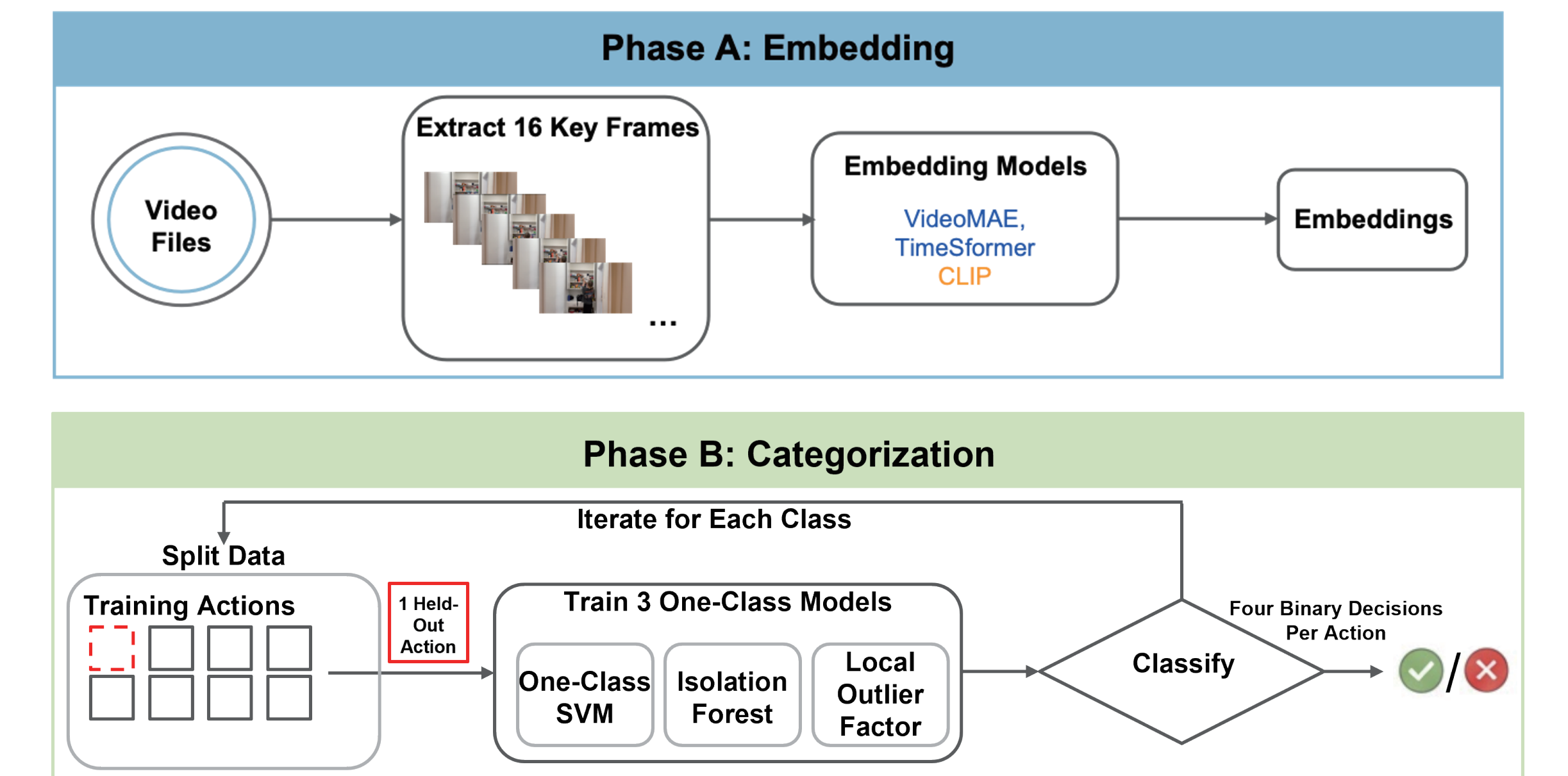
### Overall Sensitivity (d')



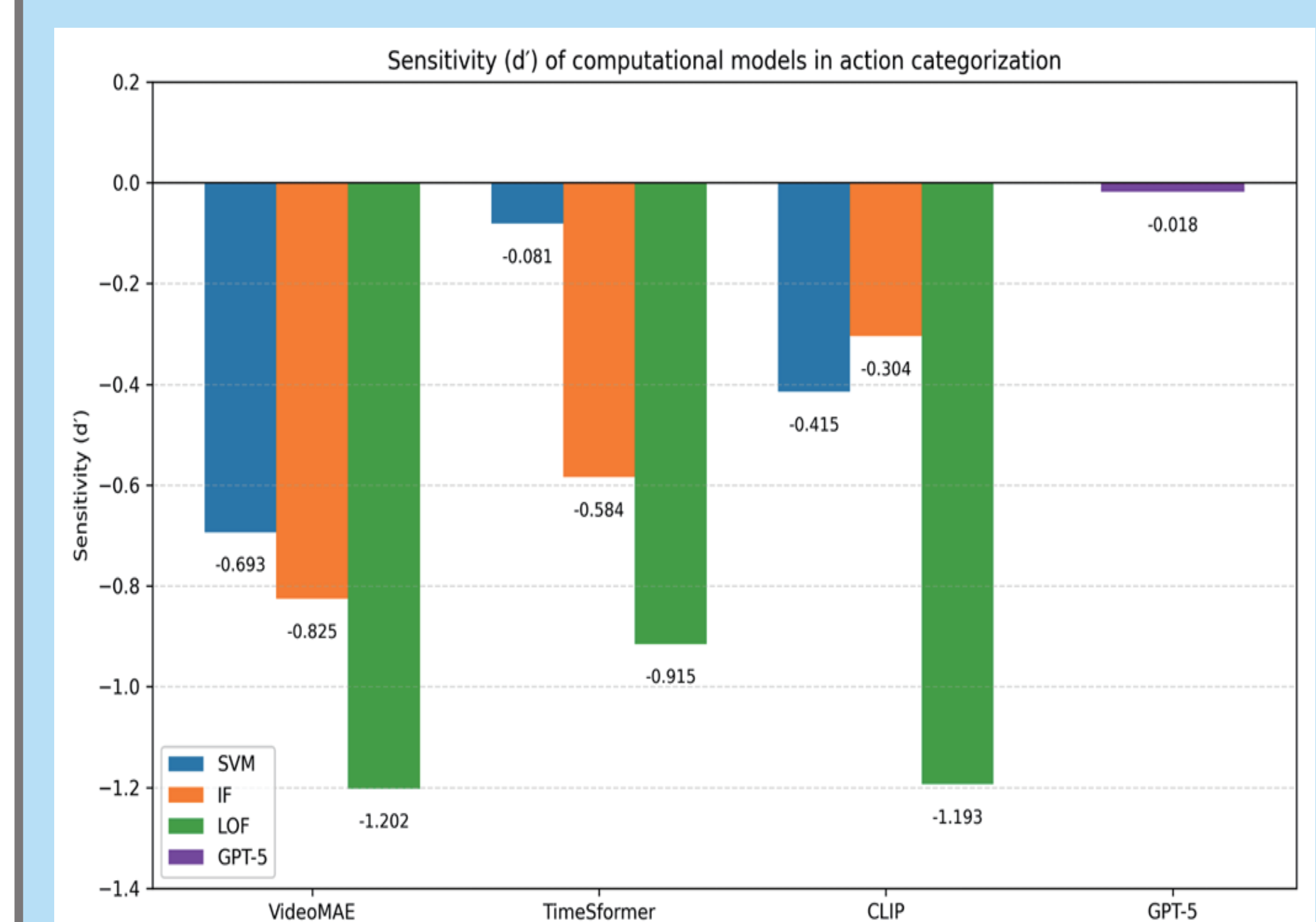
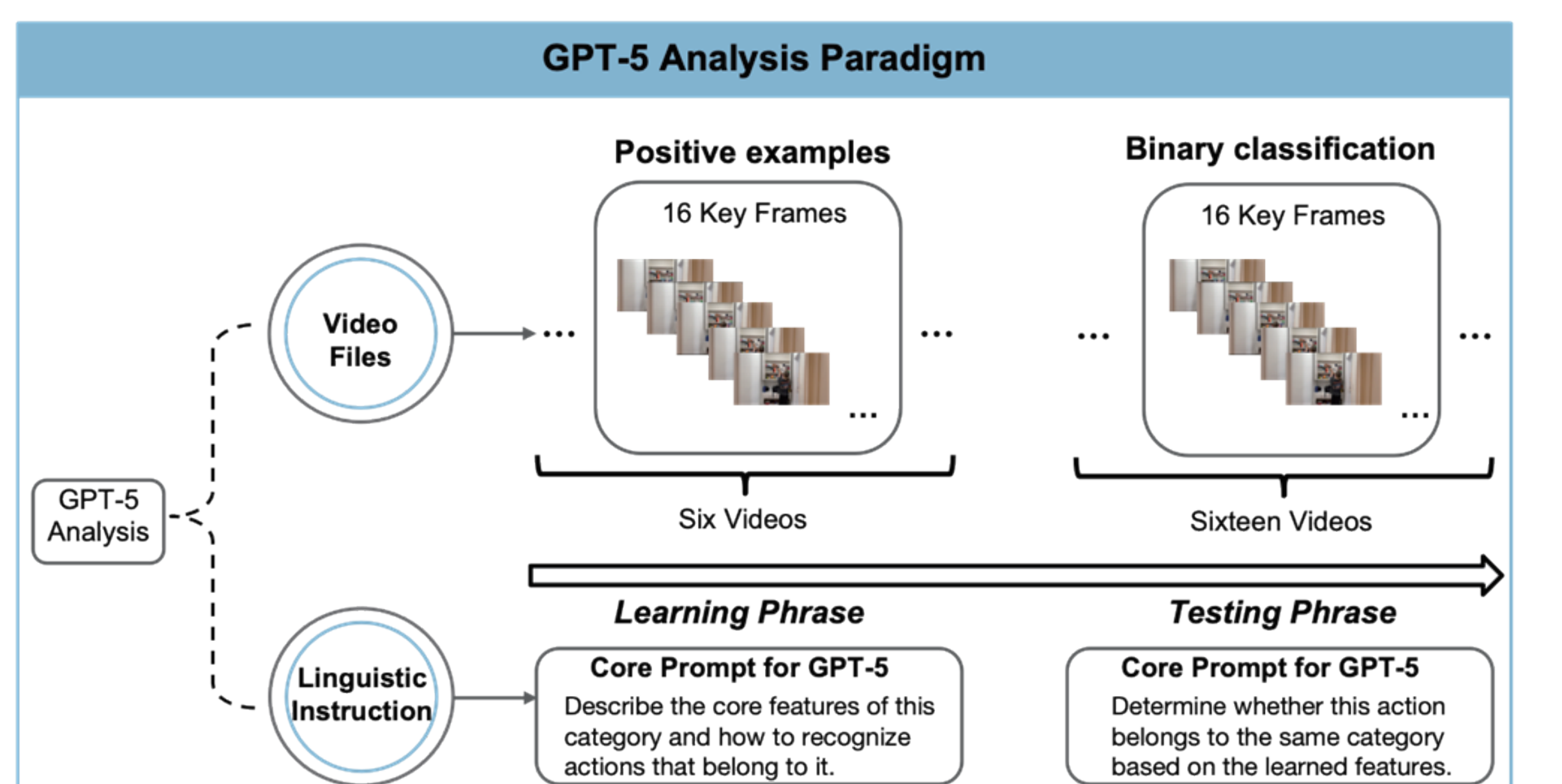
Humans reliably perform few-shot learning and generalize event categories across both language and vision; vision even supports one-shot learning.

## Expt. 4: Computational Vision Models Baseline

### Categorization Based on Models' Embeddings



### Few-shot Categorization Using GPT-5



All models show lower performance compared to chance

## Conclusion

- Humans successfully learned and generalized event categories with few-shot learning (even with one-shot learning in vision).
- Humans showed parallel patterns of event categorization in both language and vision.
- Current computational vision models lack the human-like event representation. => Suggesting that human have shared, structured event representations across language and vision, which current AI systems do not capture.