



Grounding Verbs in Intuitive Physics: Few-Shot Categorization of Action Events

Jingyi Wu^{1,2} Joshua Hartshorne³ Tao Gao²

¹Department of Linguistics, Zhejiang University

²Department of Commnuication, UCLA

³MGH Institute of Health Professions

Introduction

(No) Double-object Dative

Lafleur **slid** the puck **to** the goalie.

Lafleur **slid** the goalie the puck.

Lafleur **lifts** the crate **to** him.

Lafleur **lifts** him the crate.

Verb concepts
grounded in a
physical world model

Verb argument
structures

Physically-grounded verb argument structure

Double-object construction: Instantaneous application of force to send the object on a trajectory to a recipient

No double-object construction: Continuous application of force to an object to keep it moving

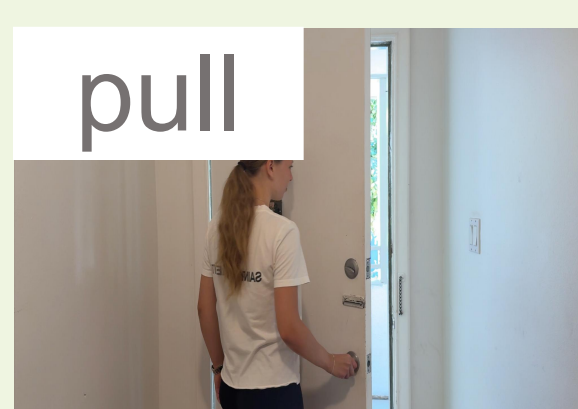
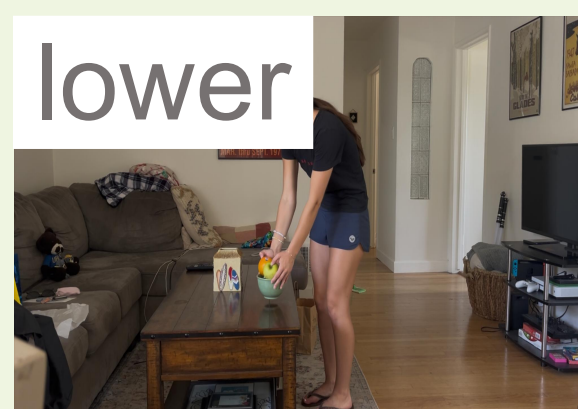
Hypothesis:

1. Humans categorize action videos based on the physical and causal attributes revealed by verb argument structure.
2. These attributes form the building blocks of human thought, hard-coded into grammar.

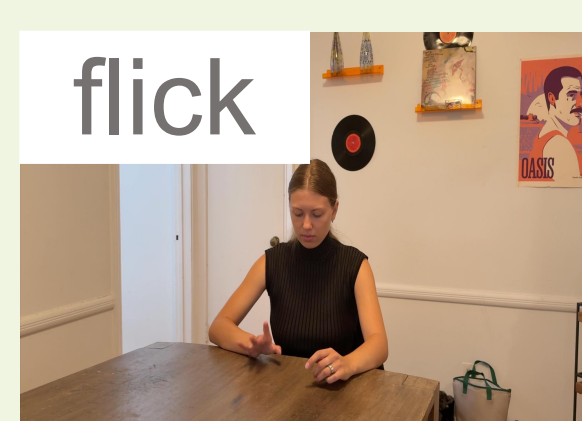
Video Dataset

4 Classes × 5 Verbs × 4 videos = 80 videos

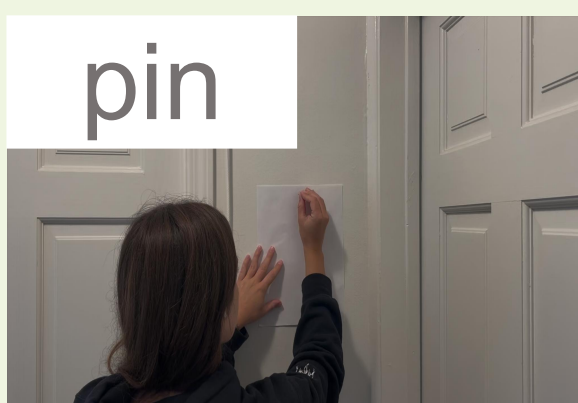
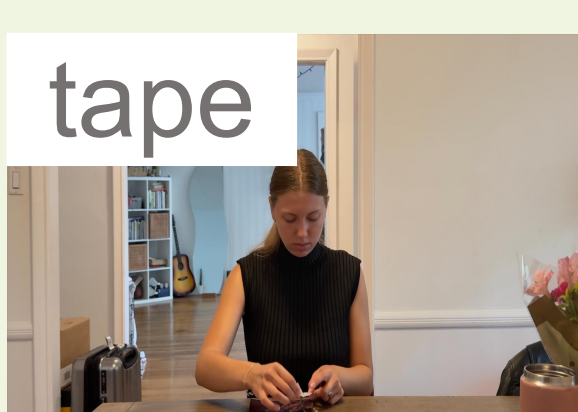
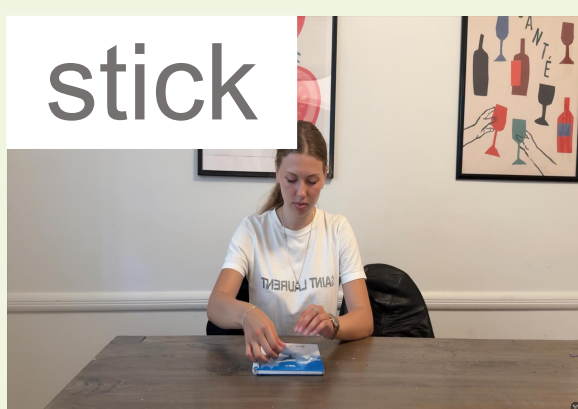
Continuous



Instantaneous



Attachment

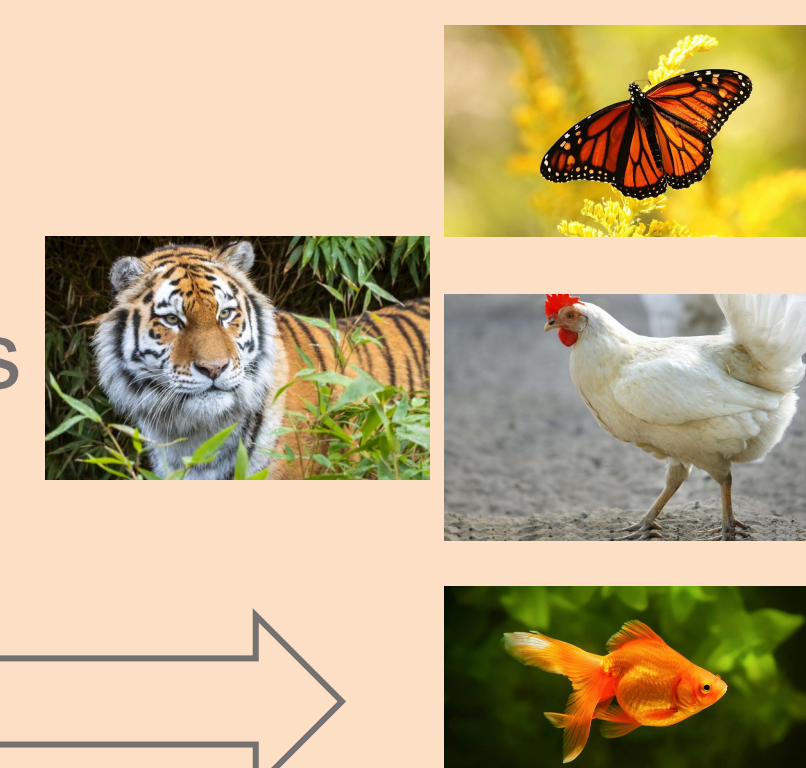


Destruction



Experiment Design

Positive: 1 mammal pictures
Negative: 3 non-mammal pictures



Practice

Learning

3 mammal animals, 2 pictures for each (3*2=6)



1 verb from each class, 1 video from each verb (4*1*1=4)

Previewing

Main Task

Learning

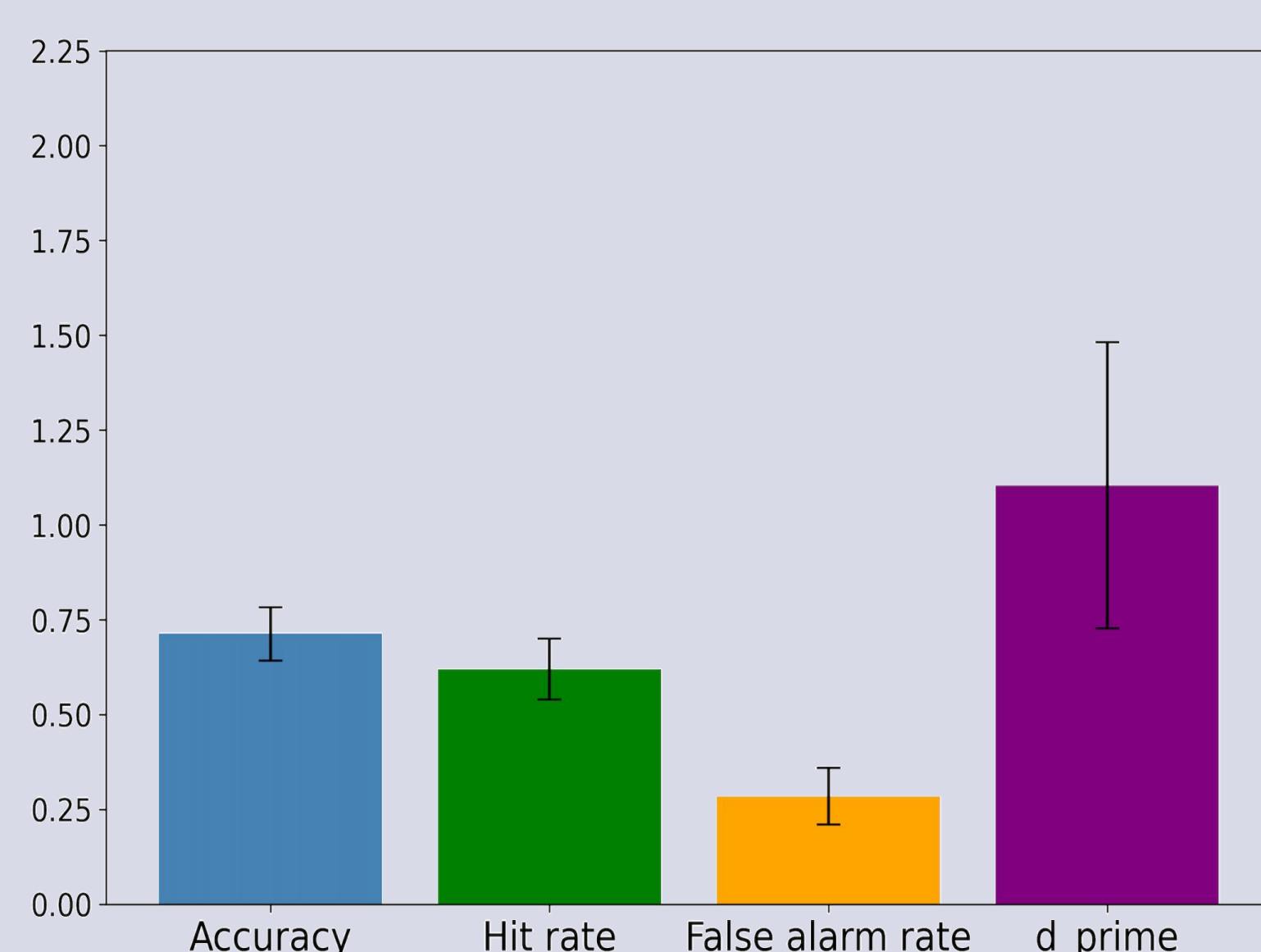
3 verbs from the testing class, 2 videos from each verb (3*2=6)

Testing

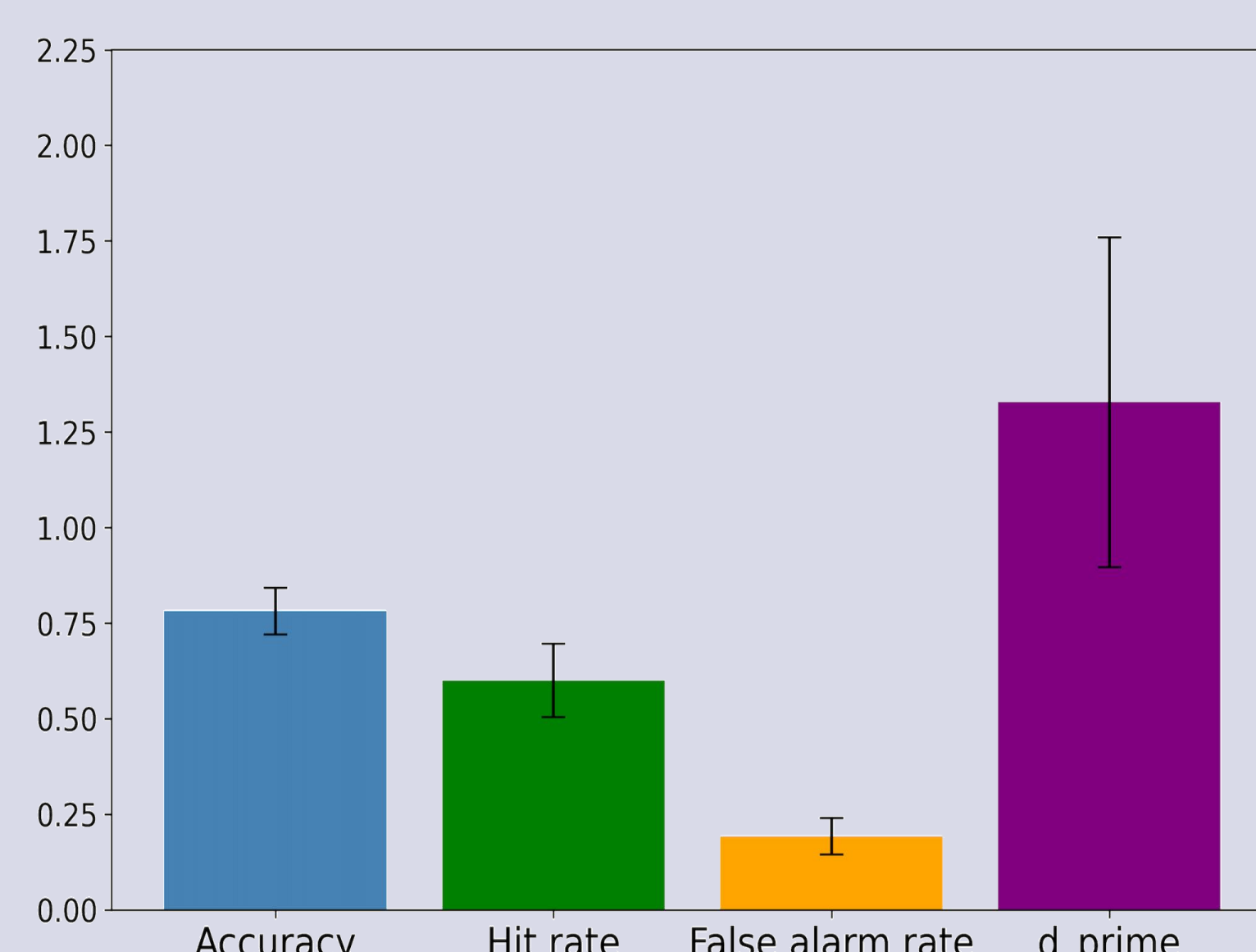
Positive: 2 unselected verbs from the testing class, 2 videos from each verb (2*2=4)
Negative: 2 verbs from each class other than the testing one, 2 videos from each verb (3*2*2=12)

Results & Discussion

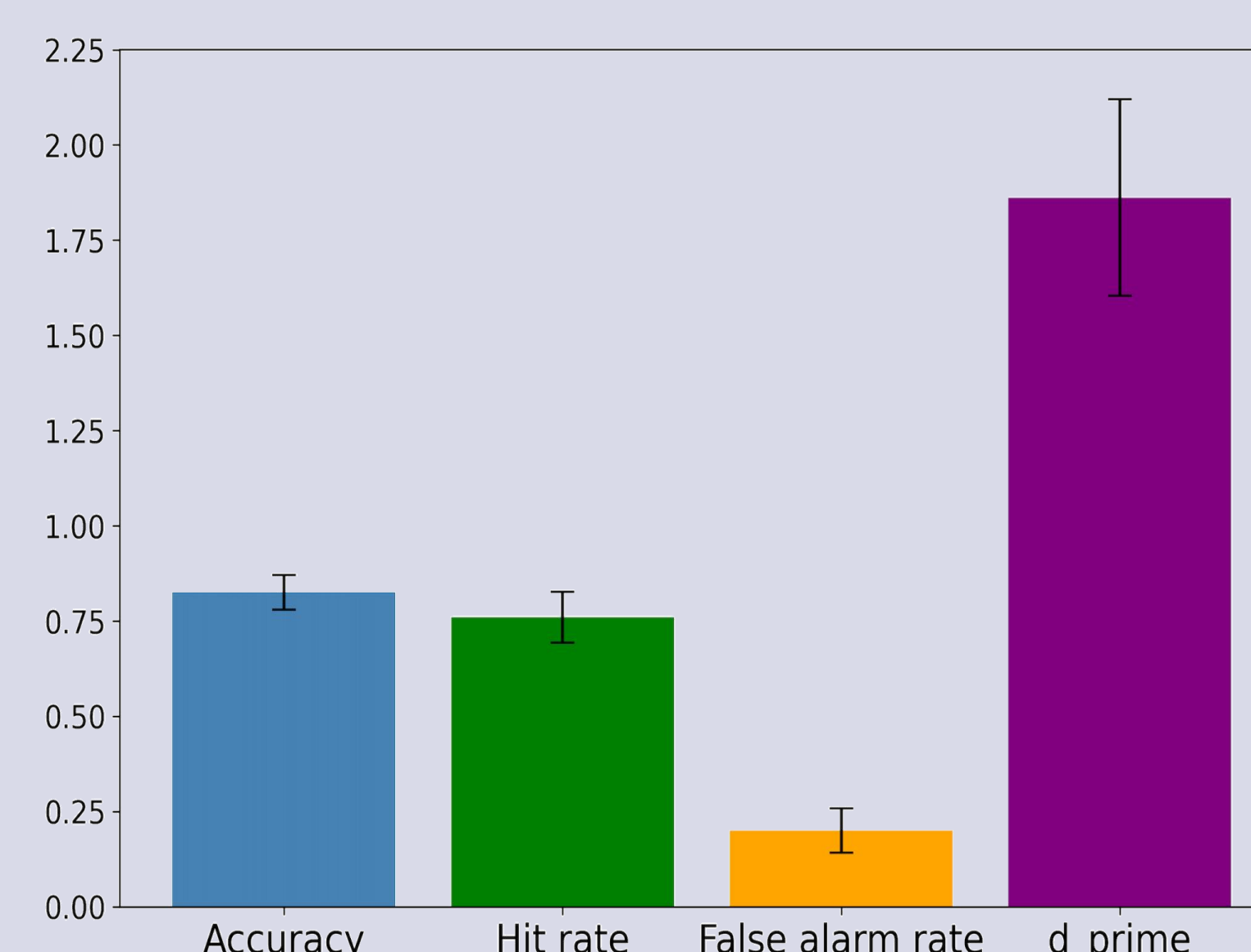
Continuous Force



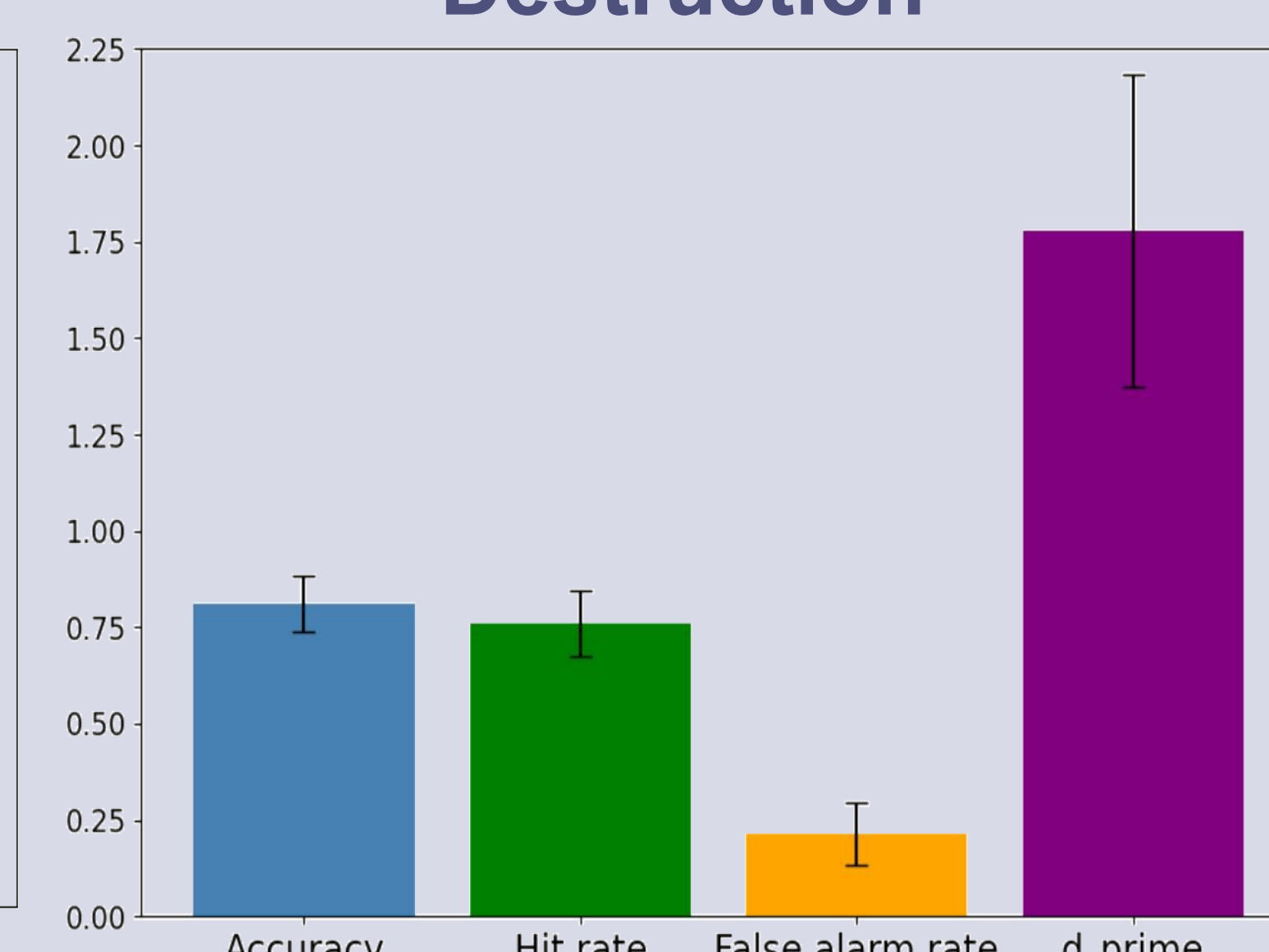
Instantaneous Force



Attachment



Destruction



Discussion:

1. **Successful few-shot event categorization:** Participants categorize unseen actions with ~70% accuracy from only a few positive examples.
2. **Variation across event categories:** Attachment and destruction are the easiest, while continuous vs. instantaneous force is often confused, with a bias to overgeneralize instantaneous as continuous.