# Experimental Contexts *Can* Facilitate Robust Semantic Property Inference in Language Models, but Inconsistently

**Kanishka Misra**[1,2]*, Allyson Ettinger[3], Kyle Mahowald[2]
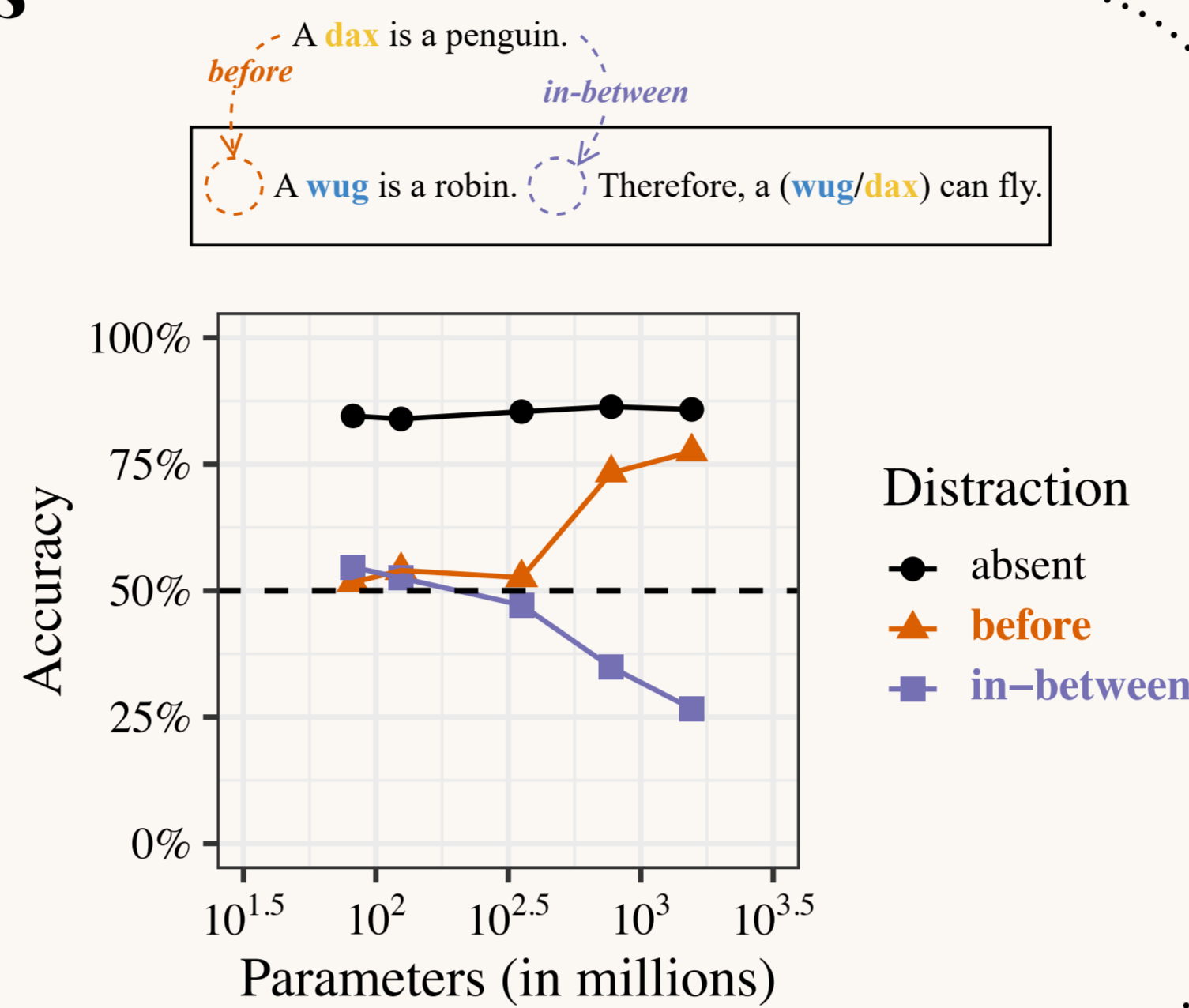
* Work done as a Postdoc at the University of Texas at Austin

[1] TOYOTA TECHNOLOGICAL INSTITUTE AT CHICAGO
[2] TEXAS — The University of Texas at Austin
[3] Ai2

SCAN ME
NSF

## COMPS: <u>C</u>onceptual <u>M</u>inimal <u>P</u>air <u>S</u>entences

A dataset to evaluate property knowledge and its robust property inheritance for novel concepts *(Misra et al., 2023, EACL)*
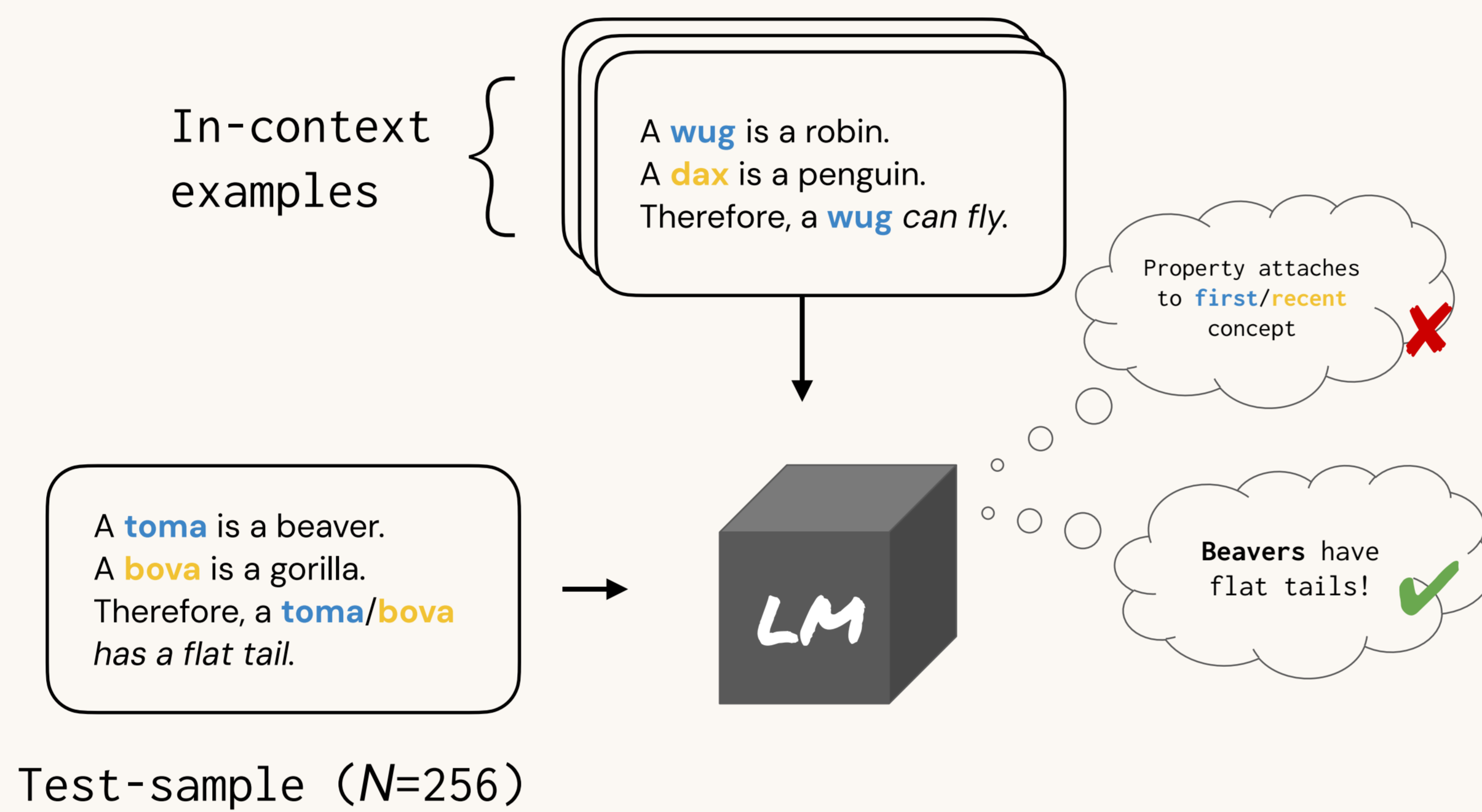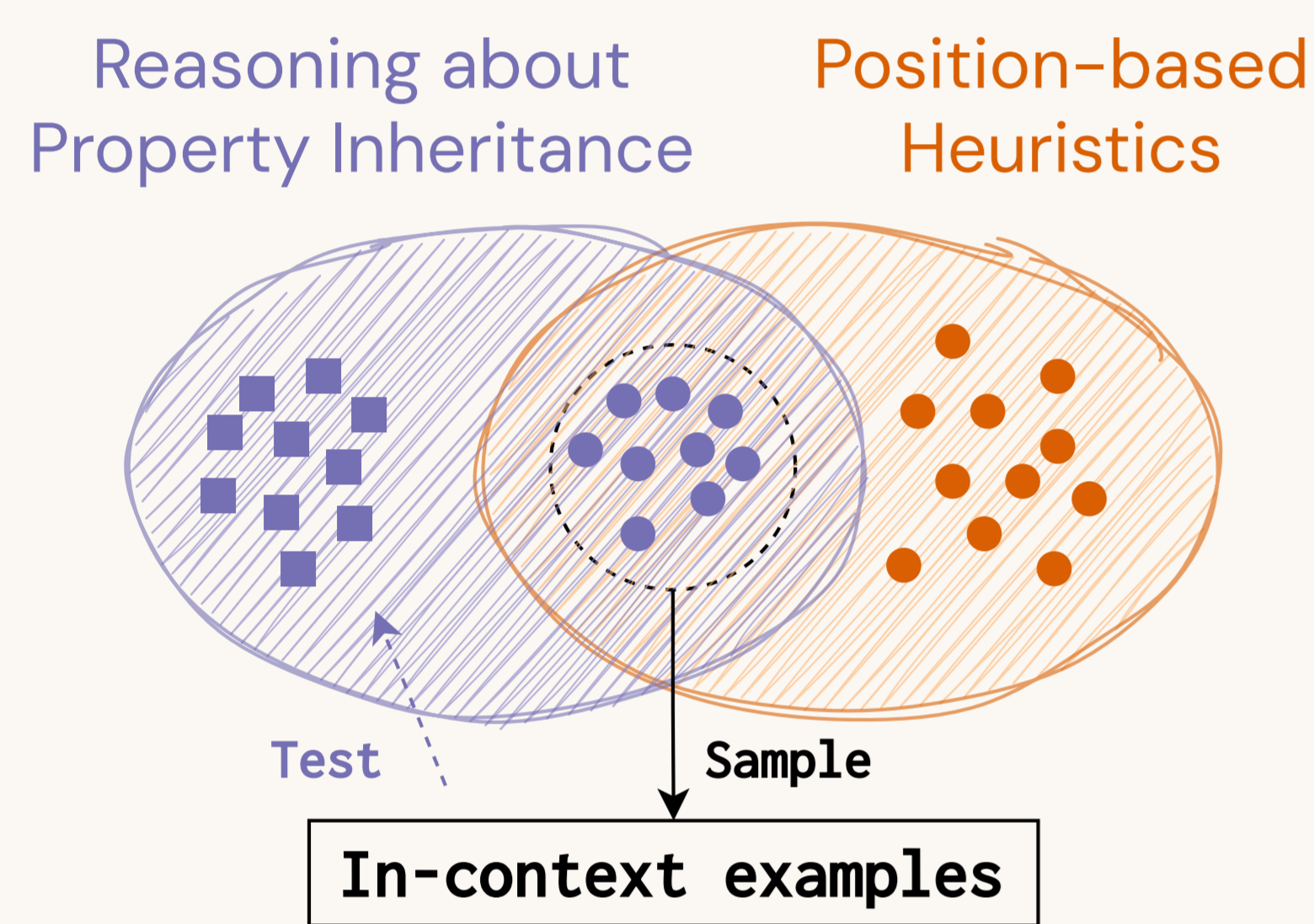
A **robin** can fly.  $>?$  A **penguin** can fly.

A **wug** is a robin.
A **dax** is a penguin.
Therefore, a **wug** can fly.  $>?$  A **wug** is a robin.
A **dax** is a penguin.
Therefore, a **dax** can fly.

A **dax** is a penguin.
*before*
*in-between*
A **wug** is a robin. Therefore, a (**wug**/**dax**) can fly.



**Premise:** LMs perform below chance when tasked to perform property inheritance for novel concepts in a zero-shot setting.

**But what happens when they are guided to an appropriate experimental context (In-context learning/instructions)?** *Lampinen (2022)*
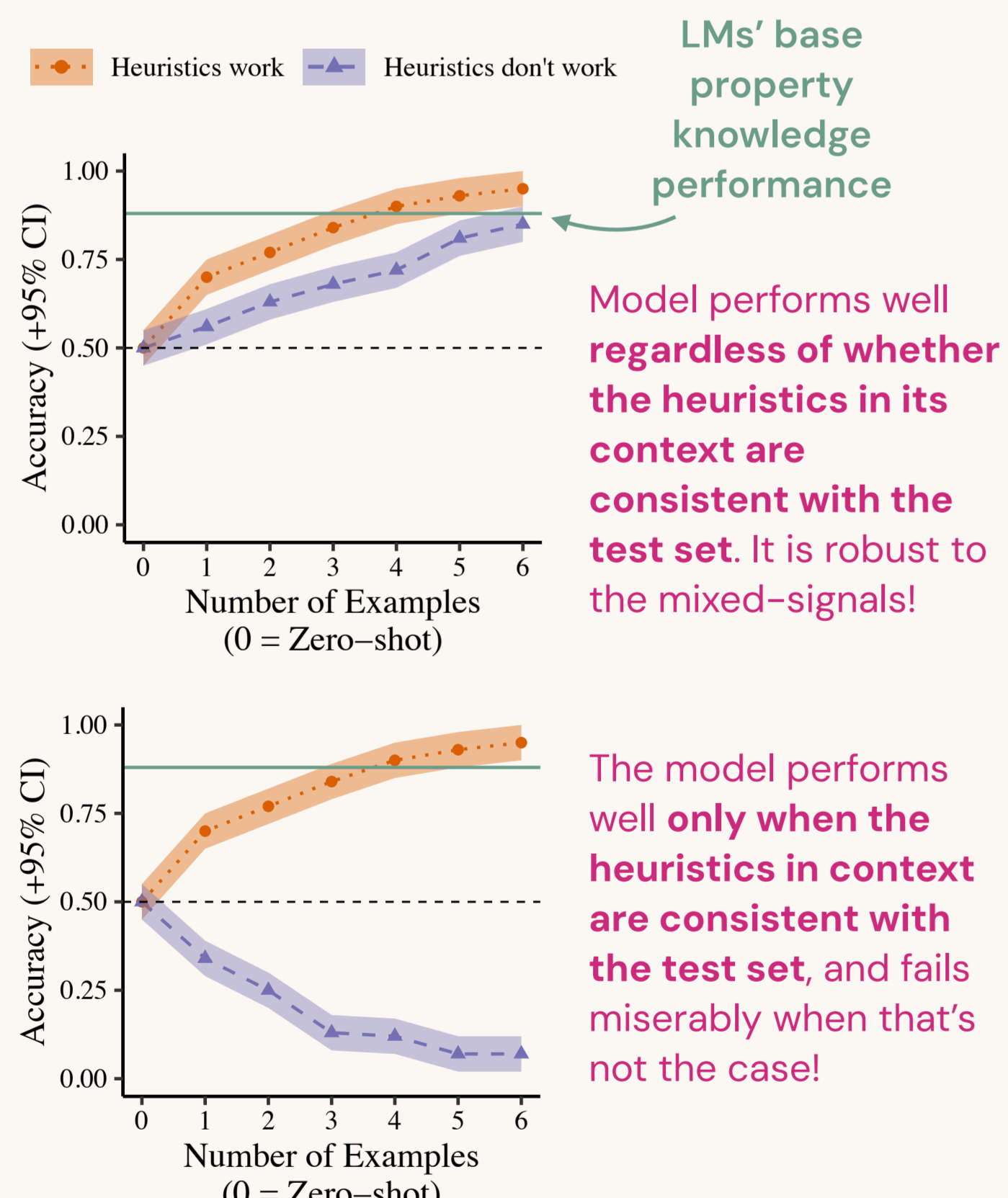
## Experiment Design

Reasoning about Property Inheritance — Position-based Heuristics

Test — Sample — **In-context examples**

In-context examples:
A **wug** is a robin.
A **dax** is a penguin.
Therefore, a **wug** can fly.

Property attaches to first/recent concept ✗

A **toma** is a beaver.
A **bova** is a gorilla.
Therefore, a **toma**/**bova** *has a flat tail*.  →  LM  →  **Beavers** have flat tails! ✓

Test-sample (*N*=256)

### Controls

- Disjointness between IC-examples and test stimuli in terms of:
  - Novel words used (wug, dax, etc.)
  - Concepts and Properties
- Multiple sets of IC examples (10) to measure variability.
- Novel words are counterbalanced (bias towards one → chance performance).
- Two types of **heuristics** tested:
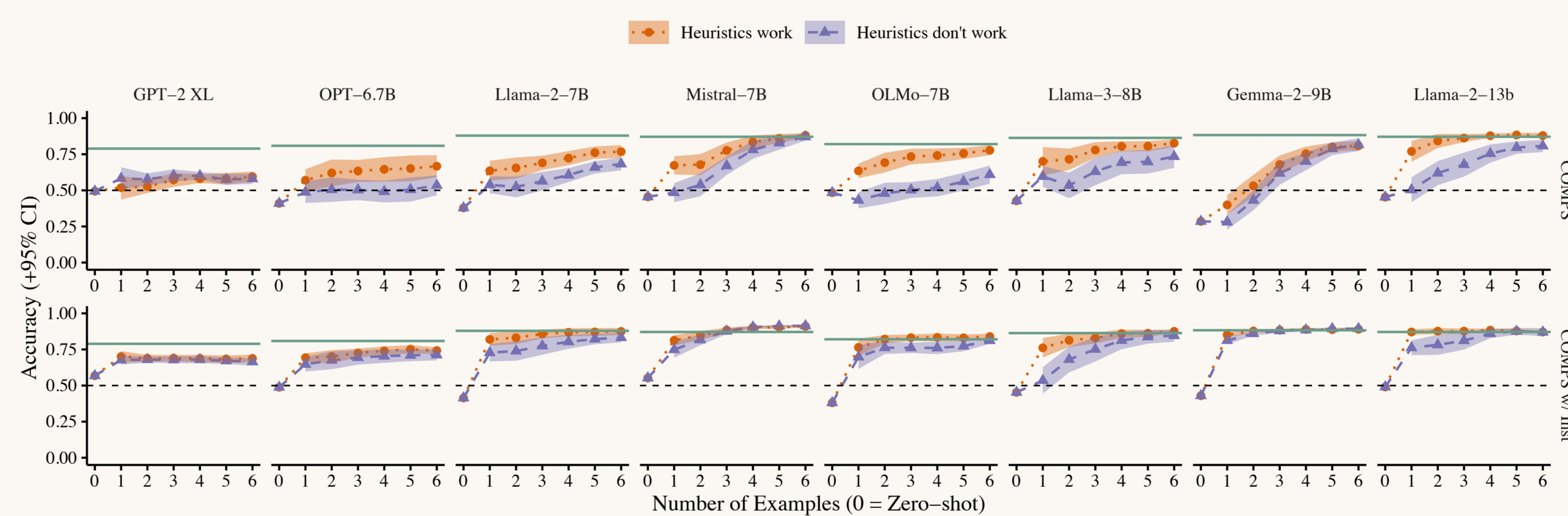  - **First correct** vs. **Recent correct**

## How to read the plots



LMs' base property knowledge performance

Model performs well **regardless of whether the heuristics in its context are consistent with the test set. It is robust to the mixed-signals!**

The model performs well **only when the heuristics in context are consistent with the test set**, and fails miserably when that's not the case!
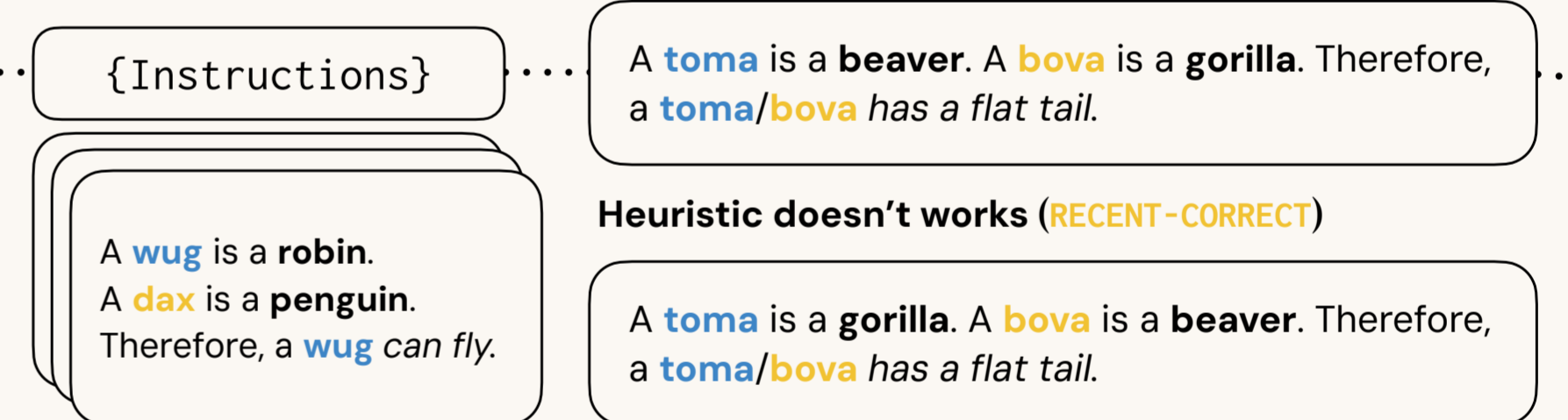
## Experimental context can improve attribution of properties to concepts...

Accuracy = Proportion of time:

$$p_\theta(\text{has a flat tail} \mid \ldots + \text{toma}) > p_\theta(\text{has a flat tail} \mid \ldots + \text{bova})$$
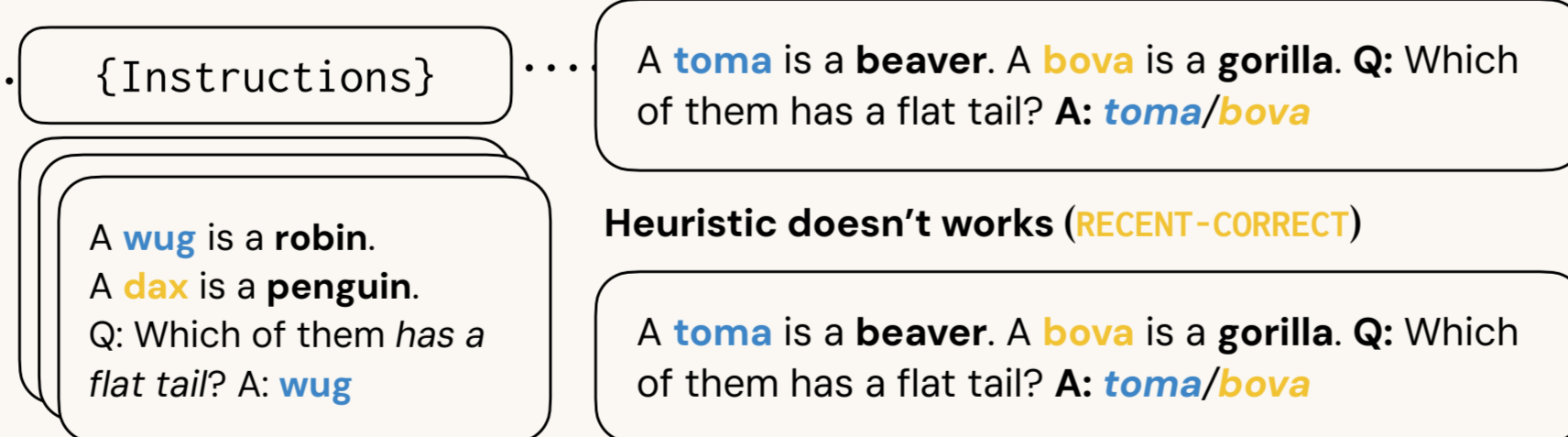


### COMPS

{Instructions}

A **wug** is a robin.
A **dax** is a penguin.
Therefore, a **wug** can fly.

**Heuristic works (FIRST-CORRECT)**
A **toma** is a **beaver**. A **bova** is a **gorilla**. Therefore, a **toma**/**bova** *has a flat tail*.

**Heuristic doesn't works (RECENT-CORRECT)**
A **toma** is a **gorilla**. A **bova** is a **beaver**. Therefore, a **toma**/**bova** *has a flat tail*.

- Experimental contexts lead to genuine improvements on COMPS.
- Instructions seem to show more robustness
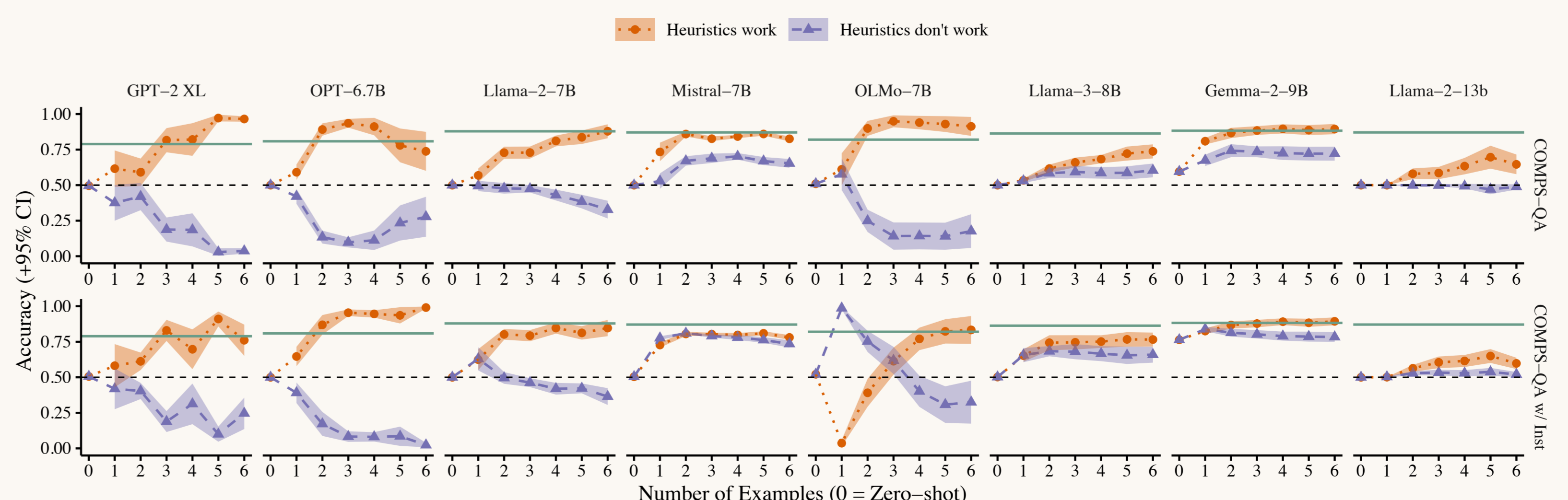- Non-trivial reliance on heuristics in some cases (OLMo, Llama-3-8B)

## ...but not the attribution of concepts to properties!

Accuracy = Proportion of time:

$$p_\theta(\text{toma} \mid \ldots + \text{Question}) > p_\theta(\text{bova} \mid \ldots + \text{Question})$$

### COMPS-QA

{Instructions}

A **wug** is a robin.
A **dax** is a penguin.
Q: Which of them *has a flat tail*? A: **wug**

**Heuristic works (FIRST-CORRECT)**
A **toma** is a **beaver**. A **bova** is a **gorilla**. Q: Which of them has a flat tail? A: **toma**/**bova**

**Heuristic doesn't works (RECENT-CORRECT)**
A **toma** is a **beaver**. A **bova** is a **beaver**. Q: Which of them has a flat tail? A: **toma**/**bova**

- Minimal reformulation of COMPS into a QA task leads to heuristic reliance in multiple models.
- **Hypothesis:** this is because the output is directly connected to the heuristic—i.e., the relative positions of concepts.