# Understanding the abilities of AI systems
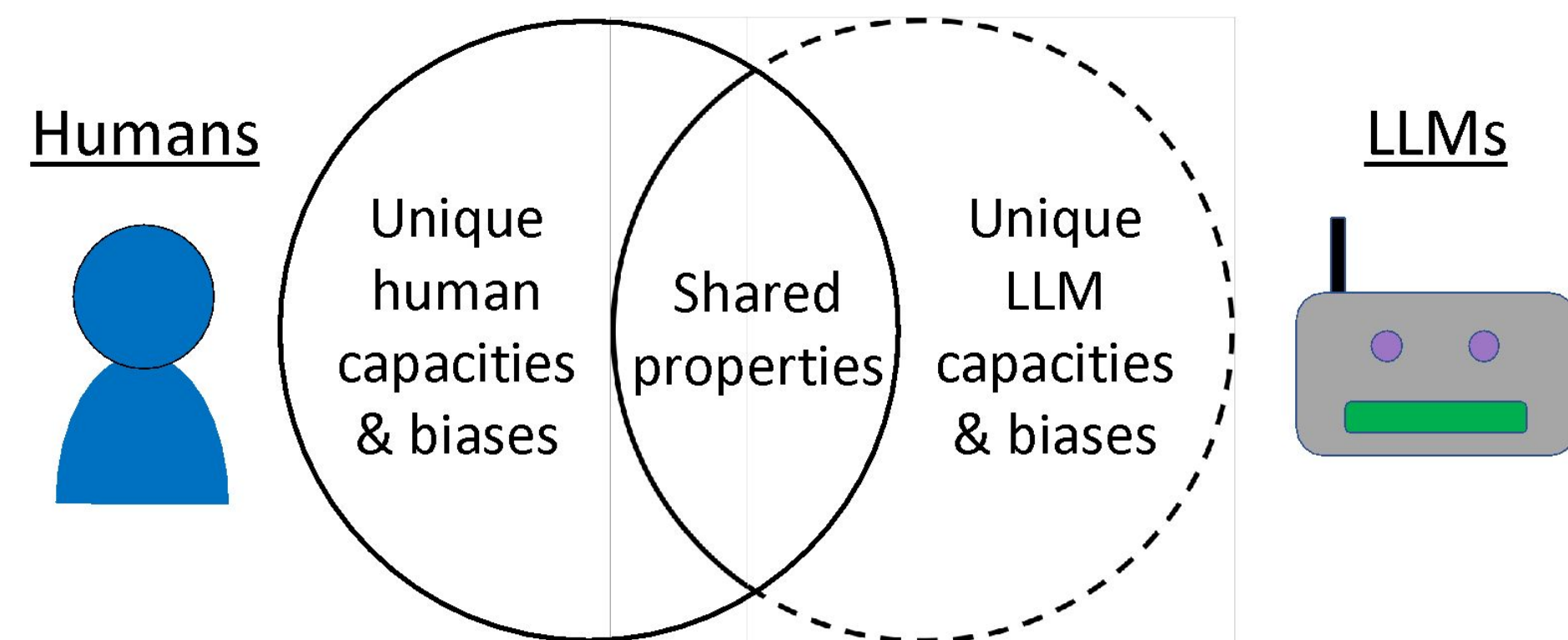
Tom McCoy,
tom.mccoy@yale.edu
Yale Department of Linguistics

Zhenghao Herbert Zhou
herbert.zhou@yale.edu
Yale Department of Linguistics

## ① Overview

- **Question:** How can we understand the (potentially non-human-like) strengths and limitations of AI systems?

- **Approach:** Analyze AI systems through the lens of the pressures that have shaped them

- **Main finding:** As predicted by our analysis, many popular AI systems are highly sensitive to probability

  - I.e., they perform better in high-probability settings than low-probability ones even when there is no difference in complexity



## ② Hypothesis: Embers of Autoregression

- Many current AI systems are large language models (LLMs)

- Primary training objective: Next-word prediction

- This objective creates pressures that favor high-probability strings of text over low-probability ones

- **Hypothesis (motivated by analyzing this objective):** LLMs will score better on high-probability examples

- All results are from Embers of Autoregression (McCoy, Yao, Friedman, Hardy, & Griffiths 2024)

## ③ Results: Output Probability

- General finding: LLMs score much better when the correct answer is a high-probability string than a low-probability one

### Example 1: Article swapping task



### Example 2: Counting letters

LLMs are much better at counting when the answer is a common number (i.e., a multiple of 10)!





## ④ Results: Task Frequency
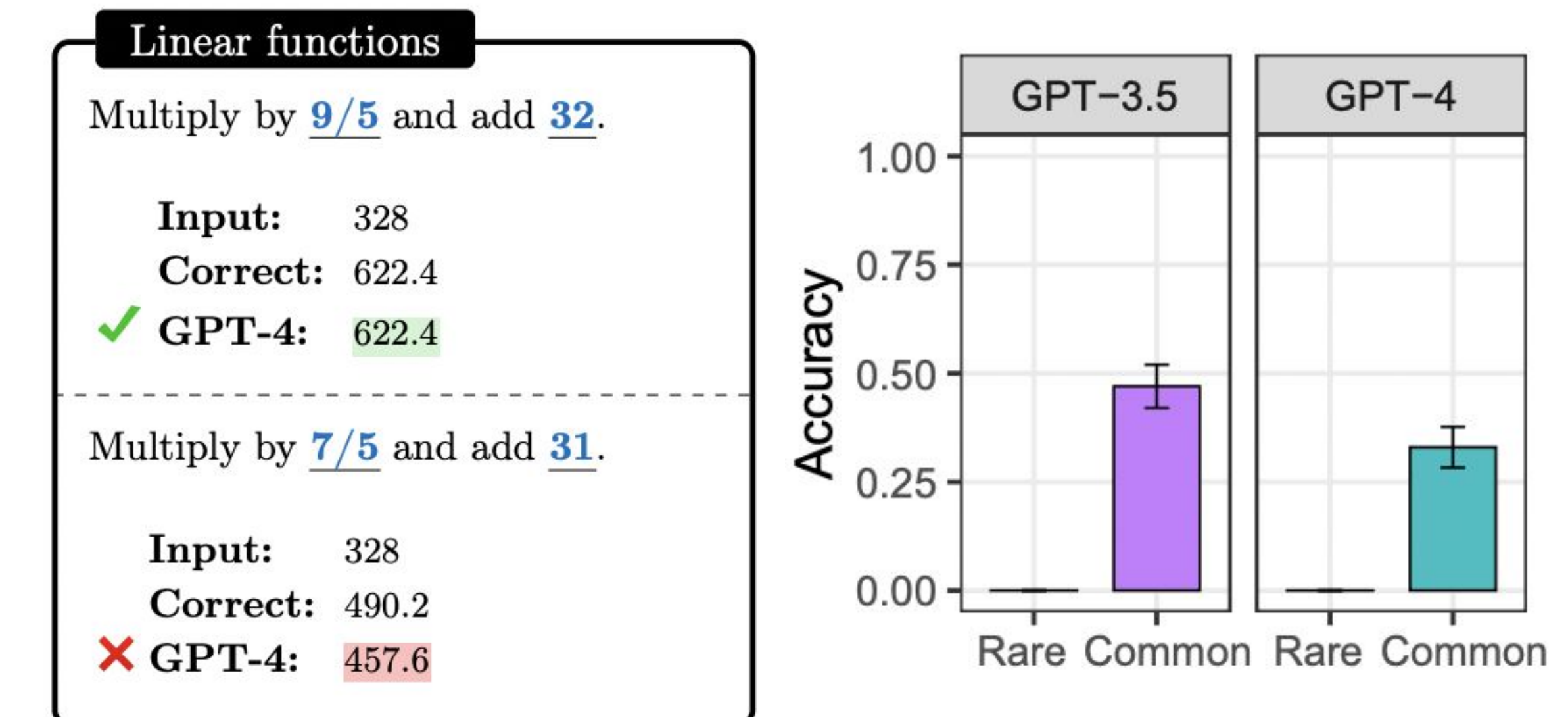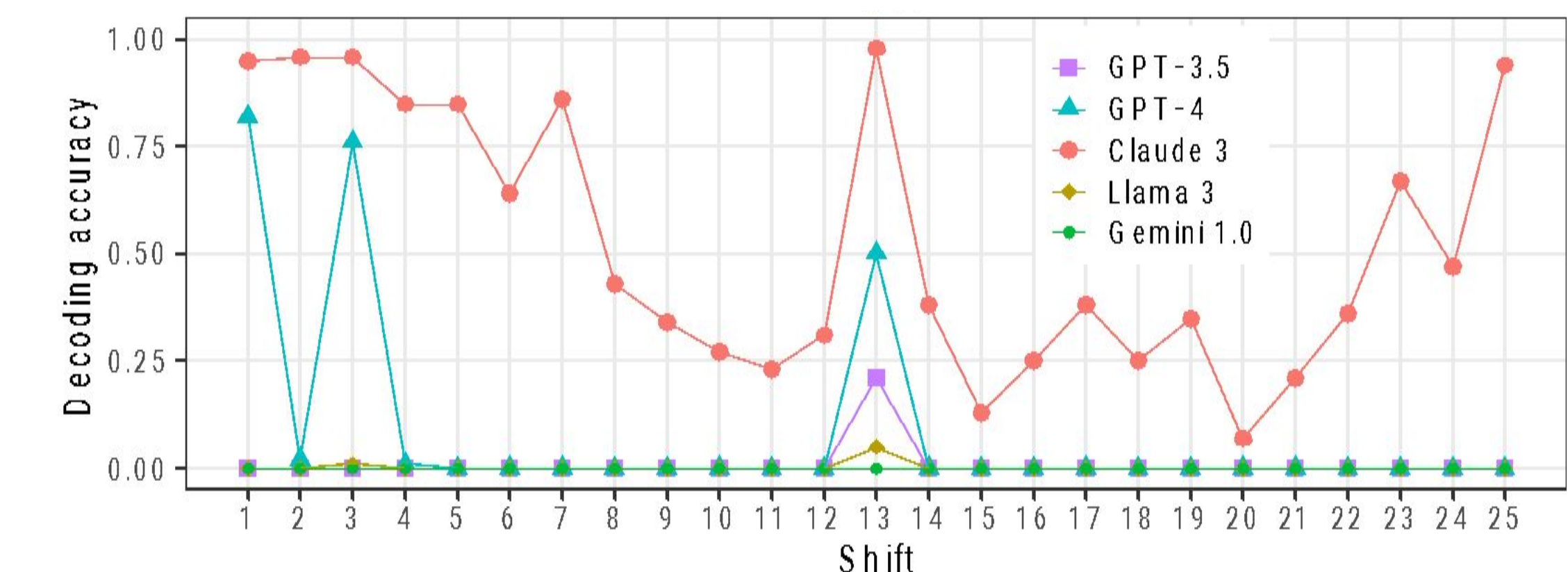
- General finding: LLMs perform much better at common task variants than rare task variants

- Example: $(9/5)x + 32$ is common (the Fahrenheit/Celsius conversion), while $(7/5)x + 31$ has no special significance



- Example: Shift ciphers are a simple type of cipher. LLMs do much better at the most common shift cipher (13) than others.



## ⑤ Conclusion

- By considering the pressures that have shaped LLMs, we predicted that they would be highly probability-sensitive

- This prediction is supported across a range of tasks

- **High-level takeaway:** To understand what AI systems are, we must understand what we have trained them to be

  - This requires thinking about the training set and assessing how the AI system does/doesn't generalize beyond it!