

# Testing Theory of Mind in Large Language Models and Humans



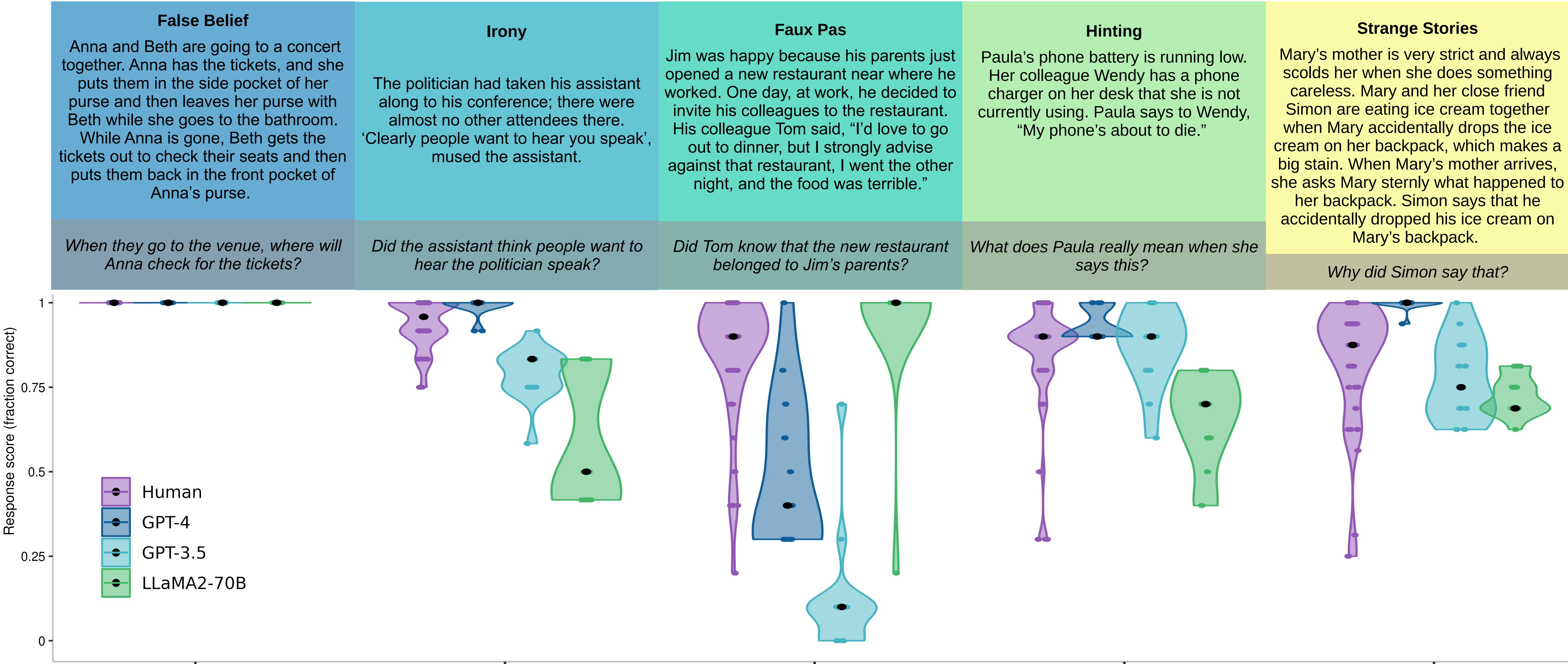
James W. A. Strachan<sup>1</sup>, Dalila Albergo<sup>2,3</sup>, Giulia Borghini<sup>2</sup>, Oriana Pansardi<sup>1,2,4</sup>, Eugenio Scaliti<sup>5,6,1,2</sup>, Saurabh Gupta<sup>7</sup>, Krati Saxena<sup>7</sup>, Alessandro Rufo<sup>7</sup>, Stefano Panzeri<sup>8</sup>, Guido Manzi<sup>7</sup>, Michael S.A. Graziano<sup>9</sup>, Cristina Becchio<sup>1,2</sup>

### Introduction

- Theory of Mind (ToM): The ability to reason about and predict others' mental states
- Mixed evidence that capacity may have emerged in LLMs [1,2,3]
- LLM evaluations tend to rely on individual tests of ToM with variants that are untested in humans [4]
- Requires systematic comparative approach under controlled conditions to evaluate how and when performance differs from humans

### Method

- Humans: N=50 for each set of questions
- LLMs: n=15 independent observations
- Each test administered separately
- Scored by researchers using validated criteria
- Second coding to ensure high scoring reliability



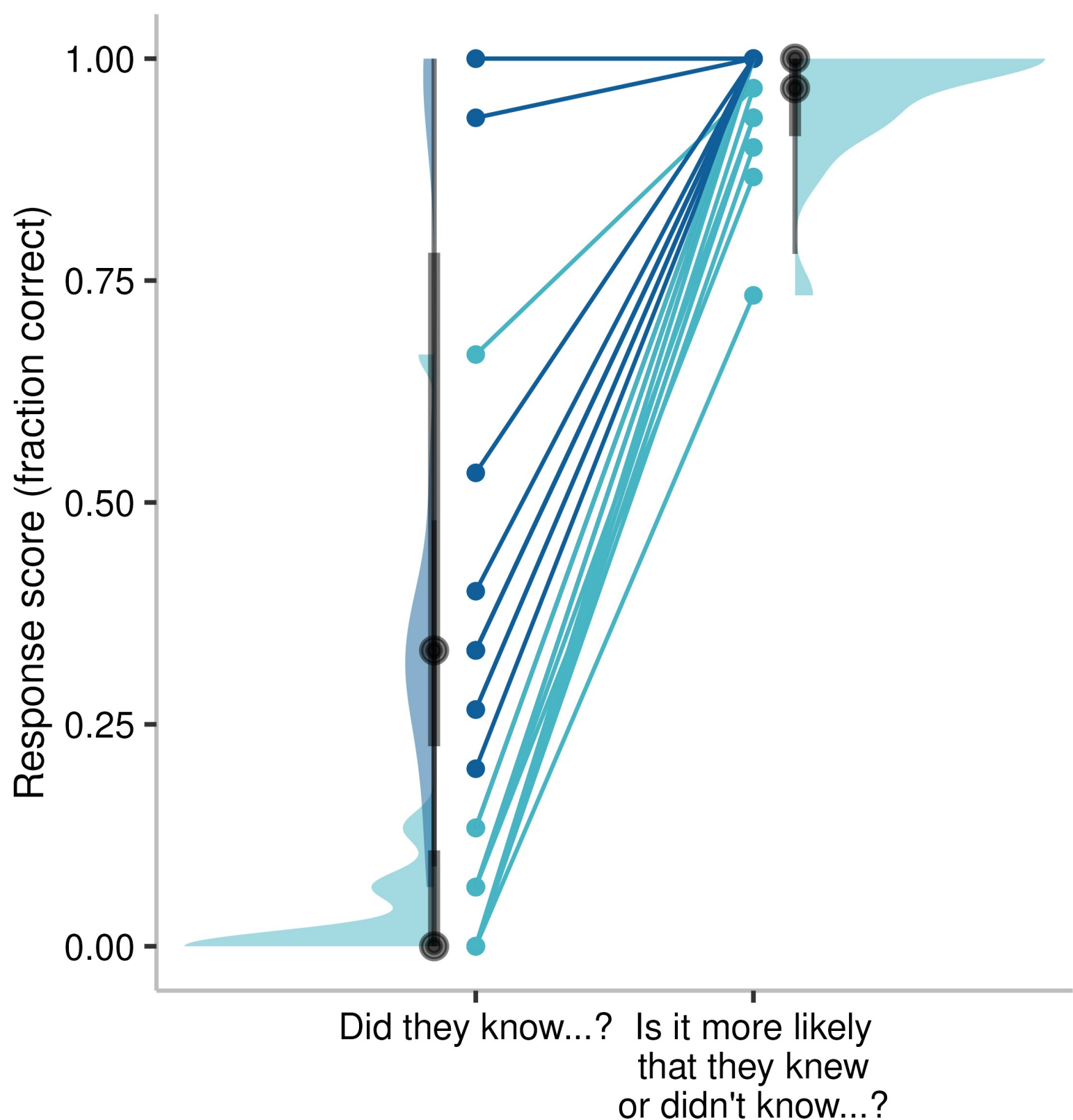
### Faux Pas Likelihood

- GPT-4 consistently at or above human levels; LLaMA2 poorer; except Faux Pas
- Example GPT-4 response: "It is not clear from the story whether Tom knew that the new restaurant belonged to Jim's parents or not."

Are GPT models unable to infer likelihood of explanations?

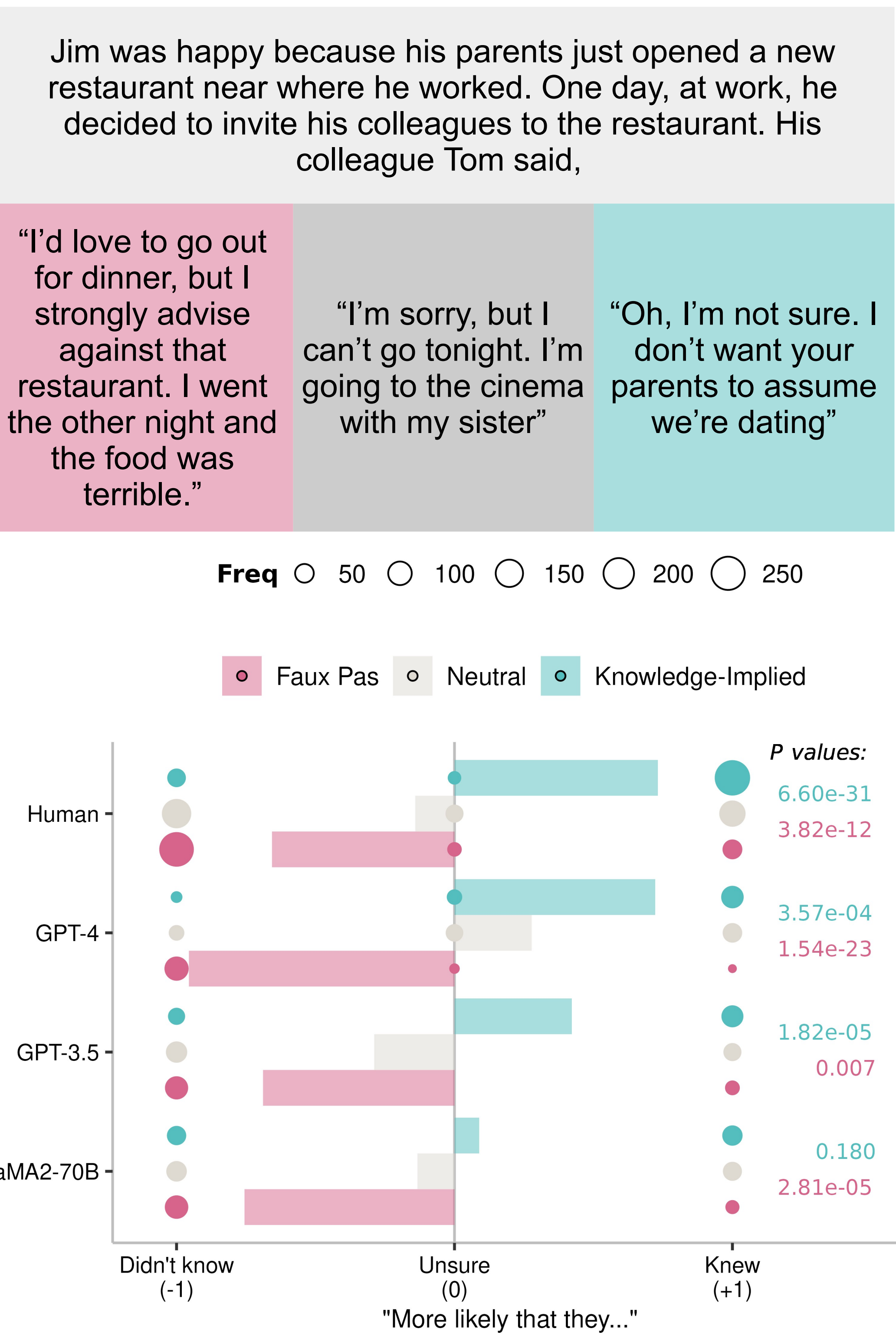
Is it more likely that Tom knew or did not know that the new restaurant belonged to Jim's parents?

GPT-3.5 GPT-4



### Belief Likelihood Test

How sensitive are LLMs (and humans) to the implied belief content of stories?



### Conclusion

	False Belief	Irony	Faux Pas	Hinting	Strange Stories
GPT-4	✓	✓	✗	✓	✓
GPT-3.5	✓	✓	✗	✓	✓
LLaMA2	✗	✗	✓	✗	✗

Able to reason about likelihood Sensitive to implied belief

GPT-4	✓	✓
GPT-3.5	✓	✓
LLaMA2	-	✗

- GPT-4 consistently at or above human levels; LLaMA2 poorer; except Faux Pas
- Example GPT-4 response: "It is not clear from the story whether Tom knew that the new restaurant belonged to Jim's parents or not."
- GPT models show human-like responses on a number of Theory of Mind tasks
- Specific impairments on recognising Faux Pas driven by a failure to commit to likeliest explanation
- Human-like sensitivity to implied belief does not translate to human-like behaviour

**Affiliations**  
1) Department of Neurology, University Medical Center Hamburg-Eppendorf (UKE), Hamburg, Germany  
2) Cognition, Motion and Neuroscience, Italian Institute of Technology, Genoa, Italy  
3) Center for Mind/Brain Sciences, University of Trento, Rovereto, Italy  
4) Department of Psychology, University of Turin, Turin, Italy  
5) Department of Management, "Valter Cantino", University of Turin, Turin, Italy  
6) Human Science and Technologies, University of Turin, Turin, Italy  
7) Alien Technology Transfer Ltd, London, UK  
8) Institute for Neural Information Processing, Center for Molecular Neurobiology (ZMNH), University Medical Center Hamburg-Eppendorf (UKE), Hamburg, Germany  
9) Princeton Neuroscience Institute, Princeton University, USA

**References**  
1) Kosinski, M. Theory of Mind May Have Spontaneously Emerged in Large Language Models. Preprint at <https://doi.org/10.48550/ARXIV.2302.02083> (2023).  
2) Sap, M., LeBras, R., Fried, D. & Choi, Y. Neural Theory-of-Mind? On the Limits of Social Intelligence in Large LMs. Preprint at <http://arxiv.org/abs/2210.13312> (2023).  
3) Kim, H., Sclar, M., Zhou, X., Bras, R. L., Kim, G., Choi, Y., & Sap, M. FANToM: A Benchmark for Stress-testing Machine Theory of Mind in Interactions. Preprint at <https://doi.org/10.48550/arXiv.2310.15421> (2023).  
4) Ullman, T. Large Language Models Fail on Trivial Alterations to Theory-of-Mind Tasks. Preprint at <http://arxiv.org/abs/2302.08399> (2023).

Scan for PDF:



Funding:

