

Metacognitive Bandits: When Do Humans Seek AI Assistance?

Aakriti Kumar, Mark Steyvers

{aakritk,msteyver}@uci.edu

University of California, Irvine

Introduction

- Humans increasingly collaborate with AI systems to make complex decisions in the real world
- But how do humans decide when to seek AI assistance?

Experimental Setup

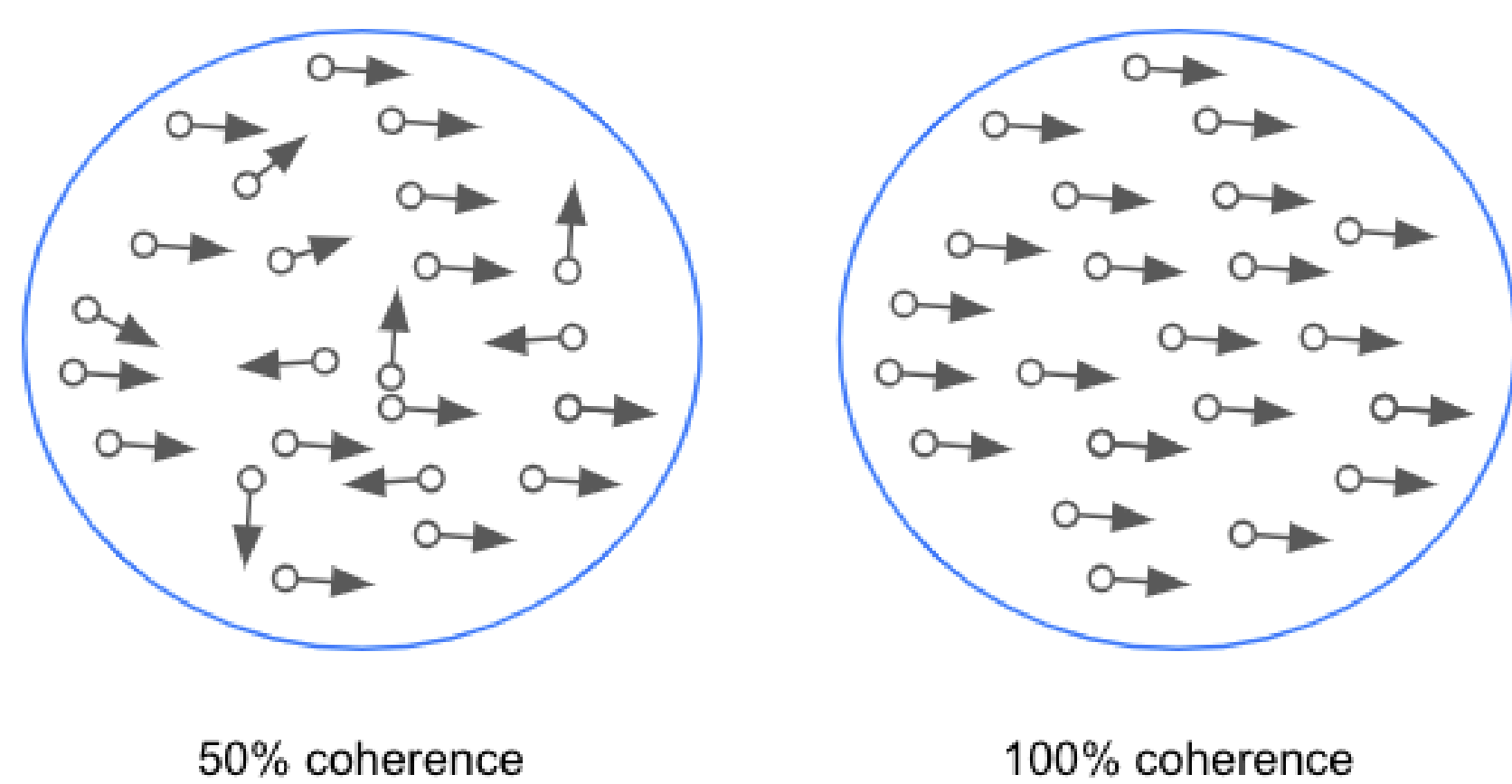


Figure: Random-Dot Kinematograms with varying coherence (inversely related to difficulty) were used as stimuli

What direction do you think the dots were moving?

Left Right

How confident are you in your response?

Low Medium High

Do you want to view the computer's decision?

Yes No

The computer thinks the dots were moving to the left

What direction do you think the dots were moving?

Left Right

You will not view the computer's decision

Continue

Correct! The dots were moving to the left

Figure: Sequence of events in the task

- Identify the dominant direction of movement in the kinematogram (left or right)
- AI advice was only shown when solicited
- AI had higher accuracy (81%) than the average participant (69%)

Metacognitive Bandits

- Decision to solicit AI advice can be explained as a combination of **explore/exploit** sequential decision making and **metacognition**
- Two armed bandit framework: decision to seek help from AI is a pull of one of two arms: Self and AI
- Arm selection depends on history of both arms and perceived difficulty of the task
- Use a Bayesian UCB framework [1] as a solution to this metacognitive task

Generative Model of Confidence & Accuracy

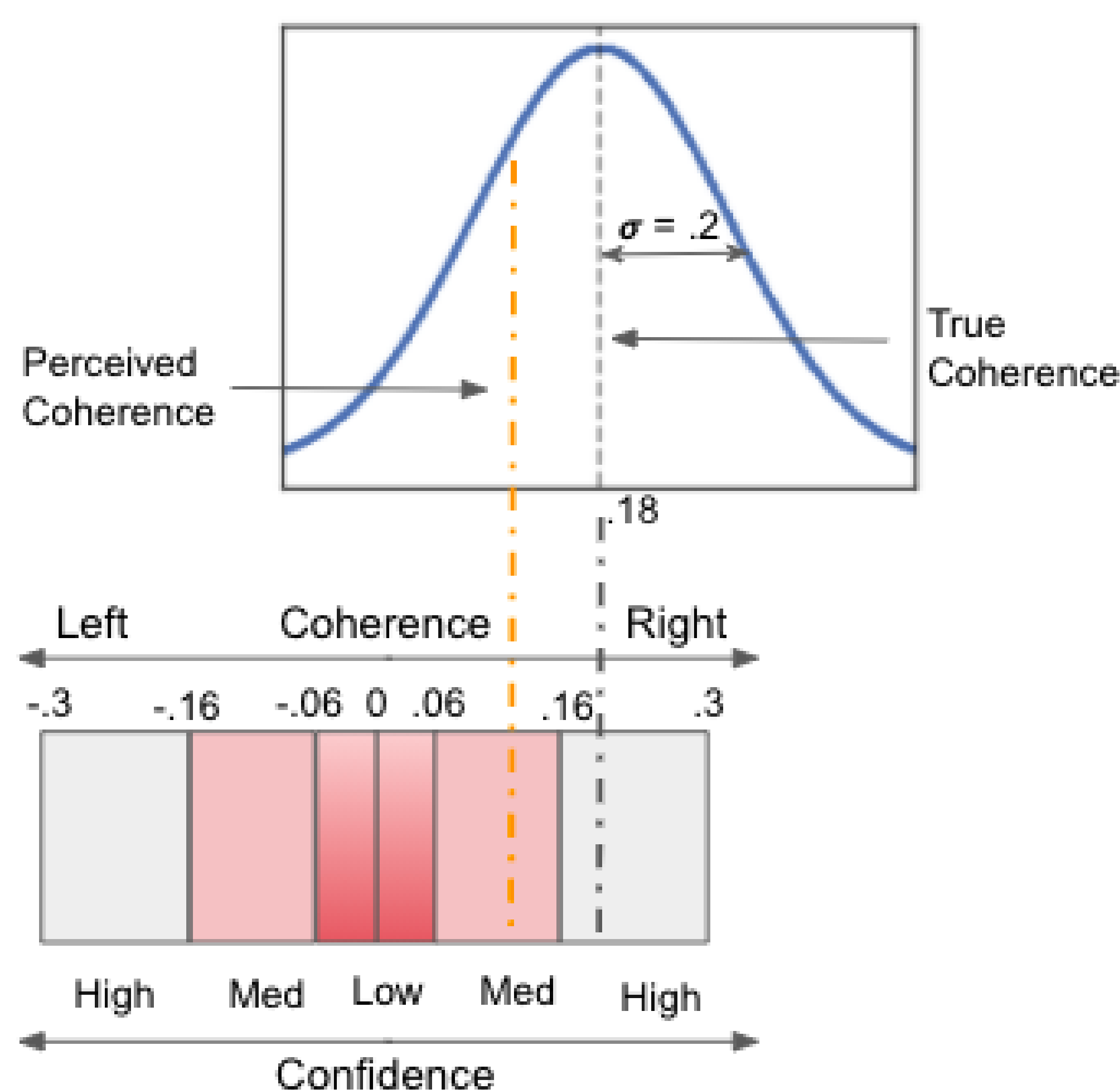


Figure: Proposed generative model for human response and confidence

- Estimated perceived coherence is used to simulate human's response and confidence on each trial
- If human's perceived coherence has the same sign as the true coherence, we predict that human gives a correct response
- Humans give higher (lower) confidence ratings for lower (higher) coherence trials

Observed & Predicted Advice Soliciting Behavior

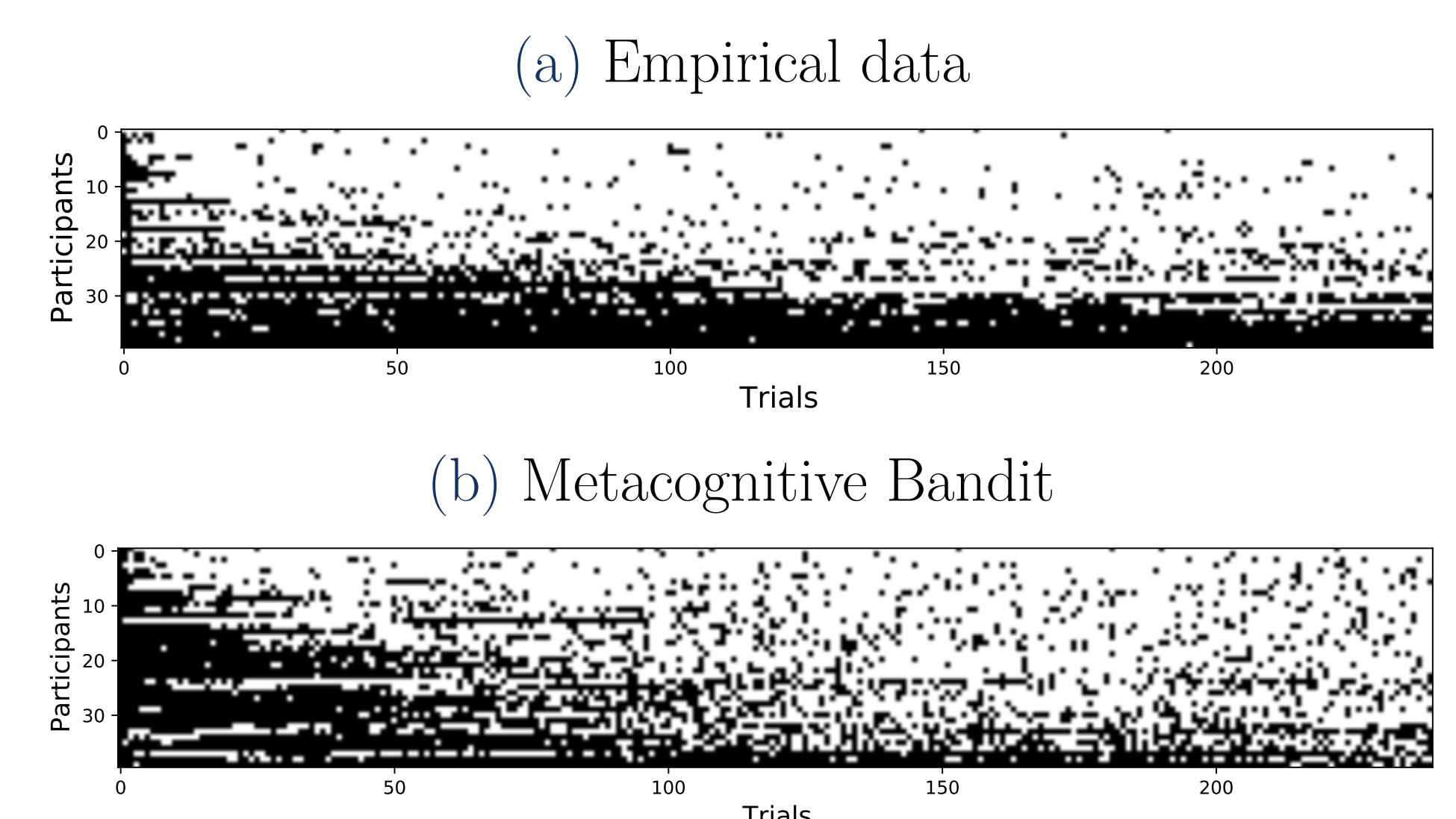


Figure: Advice soliciting behavior for actual and simulated participants on 240 trials; White corresponds to trials where a participant did not solicit AI advice.

Observed & Predicted Confidence Ratings

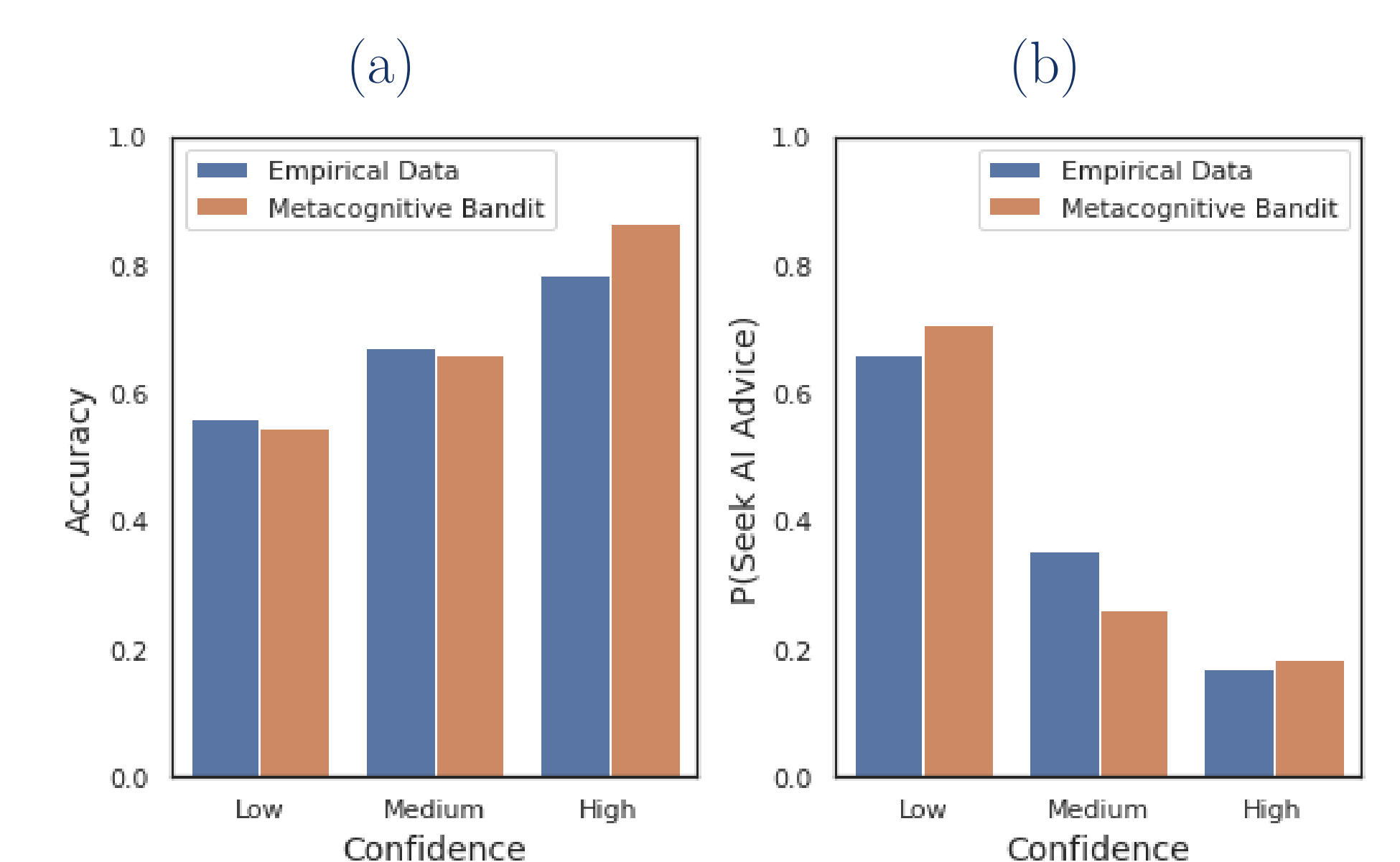


Figure: Relationship between the reported confidence of participants in their response and

- (a) the accuracy of response
- (b) probability of soliciting AI advice

Discussion & Future Work

- Currently, model only qualitatively captures trends in the data
- Look at more naturalistic decision-making settings while using a real AI in the loop
- Model how AI advice is integrated into the human's final decision

[1] Nicos G Pavlidis, Dimitris K Tasoulis, and David J Hand. Simulation studies of multi-armed bandits with covariates. In *Tenth International Conference on Computer Modeling and Simulation (uksim 2008)*, pages 493–498. IEEE, 2008.