Evaluating Risk Preferences in GPT Models

BINGHAMTON UNIVERSITY STATE UNIVERSITY OF NEW YORK

- magnitudes of potential outcomes and the associated probabilities. (1)
- human decision-making studies.
- underweight moderate- to large-probability outcomes¹. (2)
 - involves assigning disproportionate importance to it.
- lighter green indicates a higher frequency for the model to choose the safe option



Joyce Chen¹, Allen Domingo², William M. Hayes¹ Department of Psychology¹ Department of Computer Science²

BACKGROUND

• Research on human decision-making has demonstrated that individuals' risk preferences are influenced by the

• Research Objective: Measure risk preferences in GPT models and compare them to those observed in typical

• According to cumulative prospect theory (CPT), humans tend to overweight small-probability outcomes and

• Underweighting an outcome refers to assigning less importance to it, while overweighting an outcome

RESULTS

• The darker blue indicates an observed higher frequency for the model to choose the risky option, while the

• Each result is compared to its prediction under CPT's fit for each choice problem/prompt combination

METHODS



- 380 choice problems per trial, each tested 50 times
- For each trial, parameters are set such that the safe option ranges from -\$10 to +\$10, the safe option and the risky option are designed to have equal expected value.
- safe and risky options that are presented in random order with these versions of the prompt:



CONCLUSION

- Risk preferences differed between GPT-3.5-turbo and GPT-40, and preferences varied greater sensitivity to prompt format.
- All three prompts of GPT-3.5-turbo showed a bias towards the letter 'J' regardless of slightly.
- GPT-40 demonstrated high risk aversion in the gain domain for prompt A and B. The model chose the riskier options more often with loss prompts. This is consistent with prospect theory. There were minimal F/J bias.

REFERENCES

(1) Kahneman, D., & Tversky, A. (1979). Prospect Theory: An Analysis of Decision under Risk. Econometrica, 47(2), 263–291. https://doi.org/10.2307/1914185 (2) Tversky, A., Kahneman, D. Advances in prospect theory: Cumulative representation of uncertainty. J Risk Uncertainty 5, 297–323 (1992). https://doi.org/10.1007/BF00122574

DReaM Lab **Decision Research and Modeling**

while the risky option's probability ranges from 0.05 to 0.95, in increments of 0.05. Both

• The options are presented as option F and option J, which are randomly assigned to the

with different prompt templates. Compared to GPT-40, GPT-3.5-turbo demonstrated

prompt content. Prompts B and C matched the CPT models, while Prompt A differed