



# Introduction

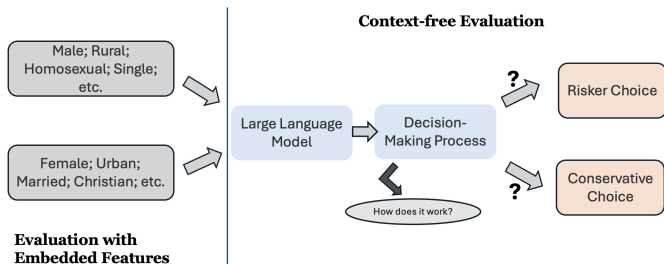
## Background and Motivation

### What is decision-making behavior?

- How agents choose between different outcomes under uncertainty.
- Key to understanding rational vs. irrational choices.

### Evaluating decision-making behaviors of LLMs:

- Increasing Use: LLMs now guide decisions in various scenerios, impacting critical outcomes.
- Need for Evaluation: Ensure LLMs make ethical and fair decisions.



# Introduction

## Background and Motivation

### Current Research Gap

- Existing evaluation models: **pre-assume** human-based norms.  
⇒ **A circular reasoning loop**: using results to validate initial questions
- **We Need:**
  - ① **A Framework**: To evaluate LLM decision-making independently of human-based assumptions.
  - ② **A Tool**: To identify fairness and sensitivity regarding various demographic features

### Research Question Statement

- **1. Evaluation Framework:**  
Assessing LLM decision-making behavior without circular reasoning logic
- **2. Fairness Issues Identification:**  
Testing both context-free and demographic-embedded scenarios

# Framework and Design

## Evaluation Framework

### Evaluation Model (utility function):

$$u(x, p; y, q) = \begin{cases} v(y) + w(p)(v(x) - v(y)) & x > y > 0 \text{ or } x < y < 0 \\ w(p)v(x) + w(q)v(y) & \text{if } x < 0 < y \end{cases}$$

where

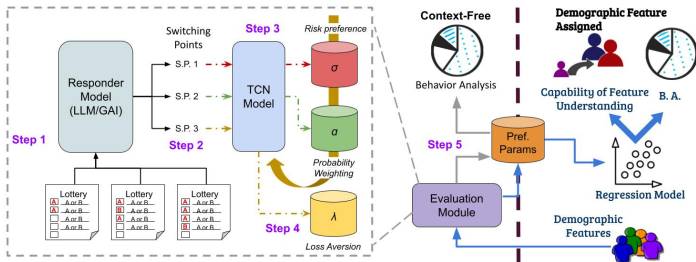
$$v(x) = \begin{cases} x^{(1-\sigma)} & \text{for } x > 0 \\ -\lambda(-x)^{(1-\sigma)} & \text{if } x < 0 \end{cases}$$

$$w(p) = \exp[-(-\ln p)^\alpha]$$

$\sigma$ : Risk Preference	$\alpha$ : Probability Weighting	$\lambda$ : Loss Aversion
<div style="display: flex; align-items: center;"> <div style="background: linear-gradient(to top, #008080, #000000, #008080); width: 20px; height: 100px; margin-right: 5px;"></div> <div style="display: flex; flex-direction: column; justify-content: space-around; width: 100px;"> <div style="text-align: center;"><b>+1</b></div> <div style="text-align: center;">↑</div> <div style="text-align: center;">0</div> <div style="text-align: center;">→</div> <div style="text-align: center;">↓</div> <div style="text-align: center;"><b>-1</b></div> </div> <div style="margin-left: 10px;"> <p>Refuse to take risks (risk-averse)</p> <p>Risk neutral</p> <p>Prefer to take risks (risk-seeking)</p> </div> </div>	<div style="display: flex; align-items: center;"> <div style="background: linear-gradient(to top, #ff4500, #ff8c00, #ffcc00); width: 20px; height: 100px; margin-right: 5px;"></div> <div style="display: flex; flex-direction: column; justify-content: space-around; width: 100px;"> <div style="text-align: center;"><b>+∞</b></div> <div style="text-align: center;">↑</div> <div style="text-align: center;">1</div> <div style="text-align: center;">→</div> <div style="text-align: center;">↓</div> <div style="text-align: center;">0</div> </div> <div style="margin-left: 10px;"> <p>Underweighting small Probabilities</p> <p>No probability distortion</p> <p>Overweighting small probabilities</p> </div> </div>	<div style="display: flex; align-items: center;"> <div style="background: linear-gradient(to top, #008000, #000000, #008000); width: 20px; height: 100px; margin-right: 5px;"></div> <div style="display: flex; flex-direction: column; justify-content: space-around; width: 100px;"> <div style="text-align: center;"><b>+∞</b></div> <div style="text-align: center;">↑</div> <div style="text-align: center;">1</div> <div style="text-align: center;">→</div> <div style="text-align: center;">↓</div> <div style="text-align: center;">0</div> </div> <div style="margin-left: 10px;"> <p>More sensitive to loss</p> <p>Neural evaluation</p> <p>More sensitive to gain</p> </div> </div>

# Framework and Design

## Experiment Setup



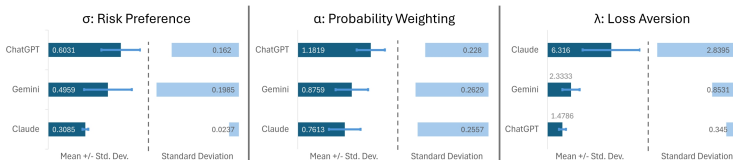
- **Step 1:** Multiple-Choice-List Experiments
- **Step 2:** Recording Switching Points

- **Step 3:** Setting Up Inequalities
- **Step 4:** Estimating Parameters
- **Step 5:** Behavior Evaluation

# Key Findings

## Basic Context-free Results

### Comparison of context-free decision-making:



### Summarization of the results:

- Each LLM model shows distinct behavior patterns, as shown in the following table:

Model	Risk Aversion	Loss Concern	Implications
ChatGPT	High	Low	Conservative, safe responses; limited novelty
Claude	Lower	High	More risk-tolerant, cautious with losses
Gemini	Balanced	Balanced	Closer to human-like behavior

# Key Findings

## Results after Embedded Demographic Features

The Personas across **10 socio-demographic groups** that we explore:

Group	Persona
<b>Panel 1: Foundational Demographic Features</b>	
Sex	male, female
Education Level	below lower secondary, lower secondary, upper secondary, short-cycle tertiary, bachelor, and graduate degrees
Marital Status	never married, married, widowed, divorced
Living Area	rural, urban
Age	15 - 24, 25 - 34, 35 - 44, 45 - 54, 55 - 64, 65+
<b>Panel 2: Advanced Demographic Features</b>	
Sex Orientation	heterosexual, homosexual, bisexual, asexual
Disability	physically-disabled, able-bodied
Race	African, Hispanic, Asian, Caucasian
Religion	Jewish, Christian, Atheist, Religious
Political Affiliation	lifelong Democrat, lifelong Republican, Barack Obama supporter, Donald Trump supporter

- We use **prompting** to embed demographic features: Assign characteristics (e.g., age, gender) to **simulate human-like decision-making contexts**.

# Key Findings

## Results after Embedded Demographic Features

### Summary:

LLMs display **different levels of sensitivity** and responses to various demographic features, influencing their decision-making behaviors.

### For Example:

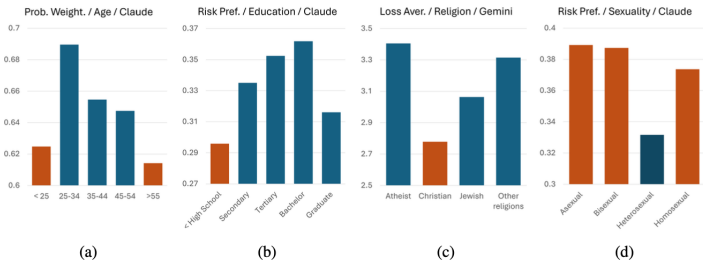


Figure: Example of Divergent Behaviors



# Implications and Discussion

## Some Open Questions

### Fundamental question:

- Should LLMs be neutral knowledge processors or reflective of human-like behaviors?

### Balancing ethical responsibility with usability:

- Should LLMs reflect human biases or aim to correct them?

### Supplementary evidence in social science research:

- Can LLM outputs help overcome survey bias in social science research, while balancing accuracy and ethical considerations?

# The End

Questions? Comments?

Please contact corresponding author:

Jingru Jia (jingruj3@illinois.edu)