

Prospect Theory Fails for LLMs: Revealing Instability of Decision-Making under Epistemic Uncertainty

Anonymous ACL submission

Abstract

Prospect Theory (PT) models the human decision-making tendency under uncertainty. While recent studies have developed some questionnaires to elicit the PT parameters (features evaluating decision-making tendency) to describe decision behavior of Large Language Models (LLM), many of them did not report the performance (or explanatory power) of PT itself for LLMs. Additionally, although PT has been used in many LLM-related fields, few studies have tried to test its robustness under linguistic uncertainty, especially epistemic markers (e.g. "maybe"). To address these research gaps, we design an experiment workflow. We adopt a classic economic questionnaire and perform parameter estimation with performance metrics (e.g. McFadden R^2). We further let LLMs make binary choices which reflect the internal probability values of epistemic markers. We then incorporate epistemic markers into the questionnaire to examine the robustness of Prospect Theory parameters. Our findings suggest that modeling LLMs' decision-making with PT is not consistently reliable, and applying Prospect Theory to LLMs is likely not robust under epistemic uncertainty.¹

1 Introduction

LLM is increasingly used in important decision-making tasks such as healthcare and finance (Keith and Stent, 2019; Lehman et al., 2022). While human decision-making under uncertainty has been extensively studied through both normative and descriptive frameworks, the decision patterns for LLMs remain underexplored. Among the human psychological model used in LLM field, Prospect Theory (PT) (Kahneman and Tversky, 1979; Tanaka et al., 2010; Rathi et al., 2025) stands out as particularly influential. It models a

¹We will release our code and data upon acceptance.

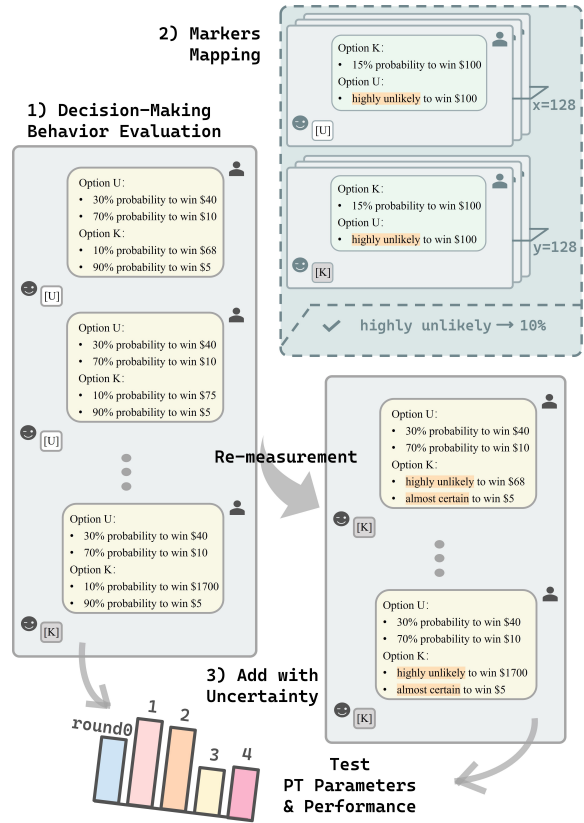


Figure 1: An overview of our three-stage experiment. Stage 1 fits PT parameters from binary choices with precise probabilities. Stage 2 estimates each markers internal probability from the point where options K and U are chosen equally. Stage 3 substitutes probabilities with markers to measure the effect of linguistic uncertainty markers.

range of human behavioral patterns by introducing psychologically grounded parameters, such as loss aversion and probability weighting (P. and H., 2005; Gneezy et al., 2006; Rathi et al., 2025). Since LLM has gained more and more attention, Prospect Theory continues to play an important role in its training, testing, and alignment (Cheng et al., 2025; Wang et al., 2025a). However, although Prospect Theory is intensively used for LLMs, few studies report the performance or explanatory power of this human psychological

052	model when applying it to LLM.	
053	In addition, decision-making is a highly safety-	
054	critical task (Leyli-abadi et al., 2025), so we must	
055	ensure the robustness of explanation theory such	
056	as PT under uncertainty. Existing studies have	
057	examined how personality prompting, sociodemo-	
058	graphic embedding, role-playing interventions, or	
059	post-training can systematically shift LLM’s per-	
060	formance (Jia et al., 2024; Liu et al., 2025a;	
061	Wang et al., 2025a), but whether and how epis-	
062	temic uncertainty markers (Lee et al., 2025) shift	
063	LLM decision-making behaviour described under	
064	Prospect Theory remains underexplored.	
065	Human uncertainty expressions are commonly	
066	expressed via epistemic markers to support flexi-	
067	ble language use (Belem et al., 2024; Liu et al.,	
068	2025b; Zhou et al., 2023). Despite their ubiq-	
069	uity, these expressions are inherently ambigu-	
070	ous (Agata, 2017), and correctly interpreting them	
071	is essential for effective decision-making (Mah-	
072	moodi et al., 2020). To examine how LLMs han-	
073	dle decision-making under uncertainty, we design	
074	a lottery-based evaluation adapted from classic	
075	behavioral economics (Allais, 1953; Brandstätter	
076	et al., 2006; Tanaka et al., 2010). Our evalua-	
077	tion consists of three stages, illustrated in Figure 1.	
078	We begin by presenting binary choice tasks with	
079	precise numerical probabilities to estimate each	
080	models Prospect Theory parameters. In the sec-	
081	ond stage, we let LLM to make binary choices be-	
082	tween a numerical probability value and a proba-	
083	bility described by epistemic markers. This allows	
084	us to infer the models implicit probabilistic inter-	
085	pretation of those uncertainty expressions. Finally,	
086	we reassess model behavior on the original deci-	
087	sion tasks, now framed using epistemic markers	
088	grounded in their inferred probability values, to di-	
089	rectly evaluate the impact of epistemic markers on	
090	decision-making.	
091	Our results reveal two key insights. First,	
092	Prospect Theory does not consistently performs	
093	well at explaining LLM decision behaviors, and	
094	it exhibits model-wise differences where larger	
095	scale models outperform the others. Second, in-	
096	troducing epistemic markers disrupts decision con-	
097	sistency and alters PT parameters, exposing the	
098	fragility of LLMs’ decision-making under linguis-	
099	tic uncertainty. These findings point to funda-	
100	mental gaps in the interpretability and robustness	
101	of current LLMs when processing ambiguous lan-	
102	guage, underscoring the need for improved epis-	
103	temic calibration and more nuanced theoretical	
	frameworks for modeling LLM behavior in uncer-	104
	tain contexts.	105
	2 Related Work	106
	Prospect Theory in Economic Decision-	107
	Making. Prospect Theory (Kahneman and	108
	Tversky, 1979) has long served as a foundational	109
	framework for modeling human decision-making	110
	under risk, capturing key behavioral patterns like	111
	loss aversion and probability distortion. Empirical	112
	studies extend PT to diverse populations and	113
	settings (Tanaka et al., 2010), and recent work	114
	explores whether LLMs exhibit similar patterns	115
	(Jia et al., 2024). Our work builds on these efforts	116
	but focuses on how language-based uncertainty	117
	impacts PT-consistent behavior in LLMs.	118
	LLM Decision-Making under Uncertainty.	119
	Recent work evaluates how LLMs handle uncer-	120
	tain scenarios, including economic games, moral	121
	dilemmas, and ambiguous instructions (Liu et al.,	122
	2025a; Jia et al., 2024; Zong et al., 2025; Zheng	123
	et al., 2025d). These studies often probe whether	124
	LLMs mimic human cognitive biases or align with	125
	normative models. We extend this line of inquiry	126
	by testing whether linguistic uncertainty conveyed	127
	by epistemic markers causes instability of LLM	128
	behaviors in decision-making situations.	129
	Epistemic Markers in NLP. Recent work inves-	130
	tigates how LLMs interpret and respond to linguis-	131
	tic signals related to uncertainty and confidence.	132
	Zhou et al. (2023) study how different epistemic	133
	markers embedded in prompts affect model pre-	134
	dictions. Belem et al. (2024) evaluate LLMs’ in-	135
	terpretation of epistemic markers, finding partial	136
	human-like behavior but systematic biases. Liu	137
	et al. (2025b) further argue that such markers are	138
	often unreliable indicators of internal confidence	139
	in LLMs. These findings suggest that models	140
	may mimic surface linguistic patterns rather than	141
	demonstrate true epistemic reasoning. Our study	142
	builds on this line of work by examining how	143
	LLMs understand epistemic markers in economic	144
	contexts with uncertainty.	145
	3 Preliminary	146
	In rational decision theory, an agent’s preference	147
	over risky prospects follows the <i>von Neumann-</i>	148
	<i>Morgenstern expected utility framework</i> . For a	149
	prospect $P = (x_1, p_1; \dots; x_n, p_n)$ yielding out-	150
	come x_i with probability p_i , the <i>expected utility</i> is	151

σ : Risk Preference	λ : Loss Aversion	γ : Probability Weighting
$+\infty$ \uparrow Risk-seeking	$+\infty$ \uparrow More sensitive to loss	$+\infty$ \uparrow Underweighting small probabilities
1 \rightarrow Risk neutral	1 \rightarrow Neural evaluation	1 \rightarrow No probability distortion
0 \downarrow Risk-averse	0 \downarrow More sensitive to gain	0 \downarrow Overweighting small probabilities

Figure 2: **Visual illustration of the three PT parameters** (σ , λ , γ). Each parameter is shown with its meaning and an interpretation of its directional significance.

computed as:

$$EU(P) = \sum_{i=1}^n p_i \cdot u(x_i), \quad (1)$$

where $u(\cdot)$ is a cardinal utility function mapping outcomes to real numbers. Human decisions are assumed to maximize $EU(P)$ under traditional Expected Utility Theory (EUT), given by von Neumann et al. (1944). However, empirical evidence systematically violates EUT assumptions. Prospect Theory (PT) addresses these anomalies through three psychological distortions of rational utility. It maintains a utility-based approach but fundamentally alters Expected Utility Theory through three key properties:

- **Risk Preference** (σ): Agents often exhibit varying degrees of risk aversion or risk-seeking behavior.
- **Loss Aversion** (λ): Losses psychologically outweigh equivalent gains.
- **probability weighting** (γ): Agents often exhibit systematic probability distortion.

To capture these characteristics, Prospect Theory introduces two specialized functions: the *value function* formalizes how the outcomes translate into subjective utility, while the *probability weighting function* captures non-linear probability perception. Together, these functions model the distorted utility calculations that characterize PT decision-making.

The **value function** $v(x)$ quantifies subjective satisfaction from outcomes relative to a reference point (zero in this study). For outcomes $x \geq 0$ (gains) and $x < 0$ (losses), the value function is:

$$v(x) = \begin{cases} x^\sigma & \text{for } x \geq 0 \\ -\lambda(-x)^\sigma & \text{for } x < 0. \end{cases} \quad (2)$$

Key parameters here are loss aversion (λ) which is a multiplier of negative perceptions, and risk preference (σ) which controls curvature (sensitivity to values).

The **probability weighting function** formalizes the transformation of objective probabilities p in subjective decision weights:

$$w(p) = \frac{p^\gamma}{(p^\gamma + (1-p)^\gamma)^{1/\gamma}}, \quad (3)$$

where γ controls the curvature of the function.

The final utility for any binary prospect in the form $P = (x, p; y, q)$ is defined as follows:

$$u(P) = \begin{cases} v(y) + w(p)(v(x) - v(y)) & [1] \\ w(p)v(x) + w(q)v(y) & [2], \end{cases} \quad (4)$$

where [1] is “if $x > y > 0$ or $x < y < 0$ ” and [2] is “if $x < 0 < y$ ”.

All functions follow the standards set in (Kahneman and Tversky, 1979).

4 Decision-Making Behavior Evaluation

Based on the above economic frameworks, we mainly adopt the three-series lottery-choice experiment developed by Tanaka et al. (2010) to get a reliable PT parameter measurement. Series 1 and 2 are designed to elicit risk preference (σ) and probability weighting (γ), while series 3 is designed for loss aversion (λ). The prospect settings are shown in Appendix B, and the prompt design is in Appendix C.

Each lottery consists of two options: a relatively safer option K with lower risk, and a riskier option U. The agent is asked to directly choose from option K and U based on its risk preference.

After sampling 256 times for each question, we count the portion of choosing option K for each lottery. Then we define the predicted probability

of choosing option K for each lottery as follows:

$$P(\text{choose K}) = \frac{e^{\text{EUDIFF}}}{1 + e^{\text{EUDIFF}}}, \quad (5)$$

where $\text{EUDIFF} = u(K) - u(U)$ is the difference in the distorted utility of prospect K and U under prospect theory. We choose the sigmoid function as it is a usual setting in economic studies, e.g., [Chakravarty and Roy \(2009\)](#). We add up the Bernoulli log-likelihood for all 35 lotteries as the negative log-likelihood function, and run MLE with this function to estimate σ , λ and γ .

To get **confidence intervals**, we use a bootstrap method ([Efron, 1979](#)) by generating simulated datasets through binomial sampling from predicted probabilities derived from the original parameter estimates. Specifically, for each observation i , we sample $\tilde{y}_i \sim \text{Binomial}(n = 1, p = \hat{p}_i)$ where \hat{p}_i is the predicted probability. The model parameter standard deviation $\sigma_{\hat{\theta}}$ is estimated from the bootstrap distribution, and the 95% confidence interval is constructed using the percentile method:

$$\text{CI}_{95\%} = \left[\hat{\theta}_{(0.025)}^*, \hat{\theta}_{(0.975)}^* \right], \quad (6)$$

where $\hat{\theta}_{(\alpha)}^*$ denotes the α -quantile of the bootstrap parameter estimates. This approach accounts for parameter uncertainty in finite samples.

We then quantify model goodness-of-fit by computing **McFadden pseudo- R^2** ([McFadden, 1977](#)), defined as:

$$R_{\text{McFadden}}^2 = 1 - \frac{\mathcal{L}_{\text{PT}}}{\mathcal{L}_{\text{null}}}, \quad (7)$$

where \mathcal{L}_{PT} is the log-likelihood of Prospect Theory and $\mathcal{L}_{\text{null}}$ represents the log-likelihood of the intercept-only model with uniform choice probabilities. This metric measures the improvement of our model over random guessing.

Finally, we calculate the **mean absolute error (MAE)** between the actual probability p_{actual} and the predicted probability p_{pred} of choosing option K derived from our PT model:

$$\text{MAE} = \frac{1}{N} \sum_{i=1}^N \left| p_{\text{actual}}^{(i)} - p_{\text{pred}}^{(i)} \right|, \quad (8)$$

where N denotes the number of observations. This provides a direct measure of prediction error.

Overall, the three PT parameters, their confidence intervals, the McFadden pseudo- R^2 , and the MAE score together describe the revealed agent’s risk attitude and our model’s explanatory power itself.

5 Decision-making Scenarios with Uncertainty

Real world decisions are often made under vague linguistic epistemic uncertainty rather than precise numerical probability. This urges for experiments on how epistemic markers influence LLM’s decision-making behaviour. We investigate how decision-making is affected when numerical probabilities are replaced by verbal probability expressions, or epistemic markers. Section 5.1 estimates their numerical equivalents via a controlled lottery experiment, and Section 5.2 applies these values in the PT framework to re-measure PT parameters. The prompt design is in Appendix C.

5.1 Probability Mapping of Epistemic Markers

Epistemic markers are inherently vague and context-sensitive ([Liu et al., 2025b](#); [Bergqvist, 2015](#)), yet they often substitute for precise numerical probabilities in practice. For LLMs, the ability to interpret epistemic markers consistently and meaningfully is critical if they are used as decision-support tools. However, there is little empirical understanding of how LLMs map epistemic markers to numerical probabilities, and whether this mapping is coherent across different models or aligned with human intuition.

No.	Epistemic Marker	Probability Mapping by Human
1	<i>almost certain</i>	95%
2	<i>highly likely</i>	90%
3	<i>very likely</i>	90%
4	<i>likely</i>	80%
5	<i>probable</i>	70%
6	<i>somewhat likely</i>	70%
7	<i>possible</i>	60%
8	<i>uncertain</i>	50%
9	<i>somewhat unlikely</i>	30%
10	<i>unlikely</i>	25%
11	<i>not likely</i>	20%
12	<i>doubtful</i>	20%
13	<i>very unlikely</i>	10%
14	<i>highly unlikely</i>	10%

Table 1: **Epistemic markers used in the experiment.** Human probability mapping comes from [Belem et al. \(2024\)](#)

To address this issue, we design a controlled lottery experiment in an economic decision-making context. Each trial presents the model with a

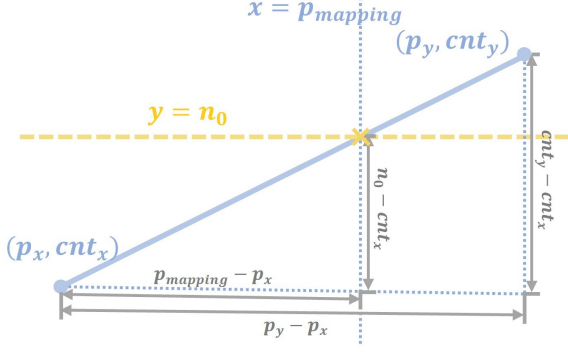


Figure 3: **An illustration of calculating $p_{mapping}$.** We use linear interpolation for fitting and calculate $p_{mapping}$ through the similarity ratio.

choice between two options. For option K, there is a fixed probability $p\%$ of winning $\$M$, while p traverses all values in set $probs$ (see Appendix A.1). For option U, there is an unknown possibility of winning $\$M$ which is defined by the different markers. For the usage of epistemic markers, we manually select 14 commonly used markers from previous work in Belem et al. (2024) to ensure that they are suitable for this context (see Table 1).

For each marker, we record the number of times the model selects option K, denoted as NUM_K . The key idea is: when NUM_K reaches $p_0 = 50\%$ of the total trials for a given probability p , the model considers the two options equally attractive. We define the mapping probability $p_{mapping}$ for that marker as the value of p at which this occurs:

$$p_{mapping} = p. \quad (9)$$

Since p is sampled at discrete points, the exact p_0 point may fall between two sampled probabilities. To estimate $p_{mapping}$, we perform linear interpolation between the two nearest points. Let $n_0 = p_0 \cdot sample_number$ be the target count corresponding to p_0 selection of option K, (p_x, cnt_x) be the probability-count pair just **below** n_0 , and (p_y, cnt_y) be the pair just **above** n_0 . As illustrated in Figure 3, by the slope formula

$$\frac{p_{mapping} - p_x}{p_y - p_x} = \frac{n_0 - cnt_x}{cnt_y - cnt_x}, \quad (10)$$

we solve for $p_{mapping}$:

$$p_{mapping} = \frac{(n_0 - cnt_x) \cdot p_y + (cnt_y - n_0) \cdot p_x}{cnt_y - cnt_x}. \quad (11)$$

Through this experiment, we obtain a list of 14 probability values (one per marker) for each

model, capturing how it semantically interprets verbal uncertainty in economic terms.

Model	σ	λ	γ	MAE↓	R^2 ↑
Human	0.670	2.630	0.685	-	-
Llama-3.1-8B-Instruct	0.585	0.010	0.753	0.332	0.092
Mistral-7B-Instruct-v0.3	0.534	0.570	0.577	0.155	0.132
Qwen2.5-7B-Instruct	0.429	0.010	3.645	0.047	0.116
Qwen2.5-14B-Instruct	0.503	1.909	0.896	0.257	0.067
Qwen2.5-32B-Instruct	0.598	1.213	0.867	0.161	0.225
gpt-5-mini	0.447	4.000	2.888	0.115	0.206
Gemini-2.5-flash	0.495	1.499	0.814	0.134	0.225

Table 2: **Baseline PT parameter estimation across LLMs.** Human reference values are taken from Tanaka et al. (2010). For *gpt-5-mini* and *Gemini-2.5-Flash*, we use *reasoning-enabled prompts* while explicitly instructing the models *not to compute expected values*. We **bold** the results where $MAE \leq 0.20$ and $R^2 \geq 0.10$, indicating that PT provides a reliable fit.

5.2 Re-measurement of Prospect Theory Parameters

After we have established the probability mapping, we pick a pair of epistemic markers which has a close normalized probability to the original settings. This is to ensure that the probability setting does not change much. Detailed replacement rules are in Appendix D. For example, for a model with “*somewhat unlikely*” mapped to $p_1 = 32\%$ and “*highly likely*” mapped to $p_2 = 68\%$, we normalize them using the formula:

$$p'_1 = \frac{p_1}{p_1 + p_2}, \quad p'_2 = \frac{p_2}{p_1 + p_2} \quad (12)$$

to ensure they add up to 100%. Then we replace the probabilities in the original decision-making behavior evaluation framework with the closest epistemic marker pairs. We re-run the PT metrics measurement test using p'_1 and p'_2 and compare the result with the original study. There are four rounds in our study. We control the option to which epistemic markers are introduced. In round 1, we only introduce markers in option K, series 1 and 2. In round 2, all 3 series of option K are with epistemic markers. In round 3, all option U are with markers. In round 4, both option K and option U of all series are with epistemic markers.

6 Results and Findings

6.1 Risk Preference under Numerical Probabilities

For decision-making under exact probability, we obtain different PT parameters compared with previous works by Jia et al. (2024) and Liu et al.

Model	Top 7 Epistemic Markers						
	almost certain	highly likely	very likely	likely	probable	somewhat likely	possible
<i>Llama-3.1-8B-Instruct</i>	87.92	56.04	58.00	41.80	44.23	36.71	36.29
<i>Mistral-7B-Instruct-v0.3</i>	96.80	67.89	63.10	57.50	87.22	48.98	52.16
<i>Qwen2.5-7B-Instruct</i>	82.71	67.00	67.06	4.78	3.44	8.93	4.51
<i>Qwen2.5-14B-Instruct</i>	91.56	55.00	54.10	42.38	26.51	32.47	38.38
<i>Qwen2.5-32B-Instruct</i>	97.50	95.08	82.82	65.00	54.74	46.08	55.00
<i>gpt-5-mini</i>	97.50	70.81	67.78	50.13	48.12	44.55	4.24
<i>Gemini-2.5-Flash</i>	97.16	61.32	52.34	45.73	32.78	4.94	3.64

Table 3: **Switching probabilities (%) for top 7 epistemic markers across models.** For *gpt-5-mini* and *Gemini-2.5-Flash*, decisions are elicited using reasoning-enabled prompts with an explicit constraint against expected-value computation.

Model	Bottom 7 Epistemic Markers						
	uncertain	somewhat unlikely	unlikely	not likely	doubtful	very unlikely	highly unlikely
<i>Llama-3.1-8B-Instruct</i>	35.49	33.71	33.49	36.91	33.65	31.30	32.83
<i>Mistral-7B-Instruct-v0.3</i>	48.08	40.15	38.27	30.93	34.03	29.47	27.88
<i>Qwen2.5-7B-Instruct</i>	35.58	27.24	19.90	27.70	25.69	18.37	19.32
<i>Qwen2.5-14B-Instruct</i>	29.03	26.45	19.10	20.82	13.04	10.94	10.52
<i>Qwen2.5-32B-Instruct</i>	2.98	21.89	3.33	3.08	3.42	2.77	2.51
<i>gpt-5-mini</i>	4.16	4.78	2.93	3.17	3.52	2.57	2.58
<i>Gemini-2.5-Flash</i>	3.25	2.86	2.51	2.50	2.52	2.51	2.50

Table 4: **Switching probabilities (%) for bottom 7 epistemic markers across models.** For *gpt-5-mini* and *Gemini-2.5-Flash*, we allow natural reasoning while explicitly prohibiting expected-value calculation.

(2025a). This is likely because of using different models and prompt design. The key parameter values we obtained are shown in Table 2, and more detail can be found in Appendix E.

Not all LLMs decision-making choices can be explained well by Prospect Theory. Following human econometric model regression standards, we consider a MAE score > 0.20 means regression result not reliable, and a McFadden pseudo- R^2 score < 0.10 means using PT may be meaningless. (McFadden, 1977) *Llama-3.1-8B-Instruct* and *Qwen2.5-14B-Instruct* have too high MAE and too low R^2 , which indicates they do not fit into Prospect Theory well.

LLMs are consistent and show human-like risk preference, but not consistent and show significantly lower loss aversion than human. All models show some extent of risk aversion σ between 0.4 and 0.6 (a bit lower than human), but they are not consistent in λ value and all significantly lower than human. Moreover, some models are more sensitive to loss ($\lambda > 1$), while others seems more sensitive to gain ($\lambda < 1$). *Llama-3.1-8B-Instruct* and *Qwen2.5-7B-Instruct* have unusual boundary λ values, which implies they have an irregular choice pattern.

Majority of LLMs show human-like probability distortion. We observe that the majority of human-like probability weighting in LLMs, which is overweighting small probabilities and underweighting large probabilities. This is the opposite result given by Jia et al. (2024).

6.2 Cross-Model Comparison of Marker Mappings

The experimental results are shown in Tables 3 and 4. We also intuitively present the probability mapping of different models for different markers in Figure 4.

Different models assign different probabilities to the same epistemic markers. Across the top 7 markers, the mapping probabilities vary significantly between models. For example, "almost certain" is mapped to over 97% in *Qwen2.5-32B-Instruct* but only less than 83% in *Qwen2.5-7B-Instruct*, while "somewhat likely" ranges from 8.9% (*Qwen2.5-7B-Instruct*) to 48.9% (*Mistral-7B-Instruct-v0.3*). This highlights that language models, even within the same family, do not share a consistent internal representation of uncertainty expressions.

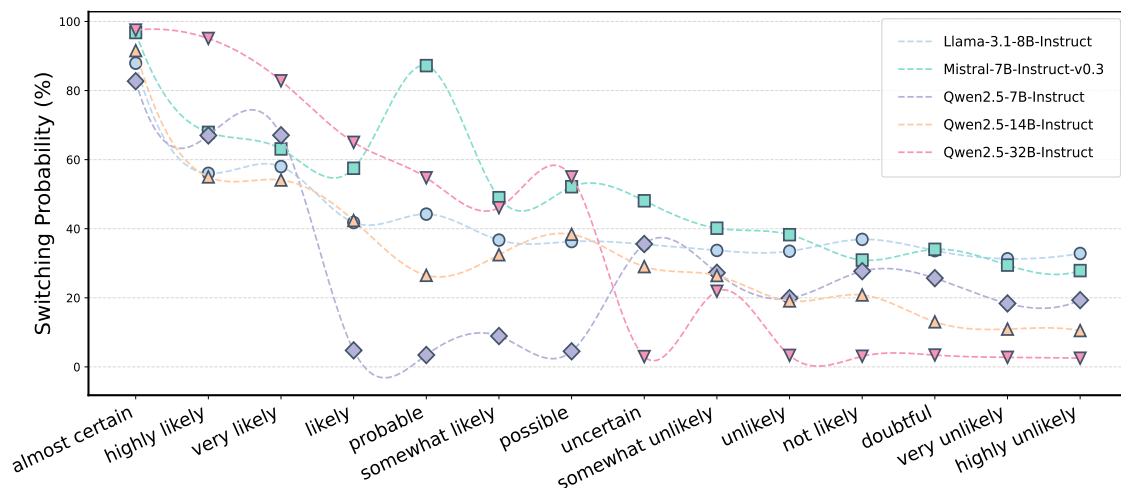


Figure 4: **Probability mapping of different models for different markers.** While different models assign divergent probability values to the same marker, the relative ordering of markers remains broadly consistent.

409 **Despite the divergence in absolute values, there**
 410 **is a striking consistency in relative ordering**
 411 **among markers.** Nearly all models assign the
 412 highest probabilities to “almost certain”, fol-
 413 lowed by “highly likely” and “very likely”, sug-
 414 gesting that models broadly agree on the ordinal
 415 semantics of epistemic language even if they differ
 416 numerically. This partial consistency could be use-
 417 ful for comparative or ranking-based tasks, but it
 418 is insufficient for applications requiring calibrated
 419 probabilistic reasoning.

420 **Some models collapse multiple distinct markers**
 421 **into indistinguishably low probabilities.** Some
 422 models exhibit a compression effect, where mul-
 423 tiple markers are mapped to similar probabili-
 424 ties. For instance, *Qwen2.5-32B-Instruct* maps
 425 6 weakest markers to probabilities between 2.5%
 426 and 3.5%, effectively collapsing fine-grained un-
 427 certainty distinctions. Such behavior suggests an
 428 overconfident or categorical interpretation of am-
 429 biguous language, which undermines the expres-
 430 siveness of epistemic markers. Similarly, *Llama-*
 431 *3.1-8B-Instruct* displays a narrow dynamic range
 432 for markers, most of them close to 30%, indicat-
 433 ing limited sensitivity of epistemic markers.

434 **Larger model size does not guarantee human-**
 435 **aligned probability interpretations.** Interest-
 436 ingly, larger models do not always produce more
 437 calibrated mappings. Although *Qwen2.5-32B-*
 438 *Instruct* maps a wider range, with some close
 439 to 100% and some near 0%, this does not pre-
 440 vent it from having poor discrimination on low-
 441 probability markers. Meanwhile, smaller mod-

442 els such as *Mistral-7B-Instruct-v0.3* perform no
 443 worse than *Qwen2.5-32B-Instruct* across the entire
 444 mapping of markers. These patterns suggest that
 445 scale alone does not guarantee a nuanced proba-
 446 bilistic understanding of epistemic markers, and
 447 that other factors such as model architecture may
 448 play a more decisive role.

6.3 Changes of Decision-Making Behavior under Epistemic Markers

449 The estimation results under varying degrees of
 450 linguistic uncertainty are illustrated in Figure 5
 451 and 6. We examine how replacing precise numer-
 452 ical probabilities with epistemic markers af-
 453 fects LLMs’ internal decision-making parameters
 454 and performance under the framework of Prospect
 455 Theory. The analysis covers four experimental
 456 conditions, including marker substitution in either
 457 or both options of the lottery choice task.
 458
 459

460 **Linguistic uncertainty causes moderate shifts**
 461 **in risk preference but affects other decision pa-**
 462 **rameters more profoundly.** For the majority of
 463 models, the estimated risk preference remains rel-
 464 atively stable despite replacing numeric probabi-
 465 lities with epistemic markers, indicating that risk
 466 attitudes are only mildly perturbed by uncertainty.
 467 However, models differ considerably in their loss
 468 aversion and probability weighting under these
 469 conditions. Larger models tend to maintain more
 470 consistent loss sensitivity resembling human be-
 471 havior, whereas smaller models show larger fluctu-
 472 ations. Similarly, the increase in γ after adding dis-
 473 turbance shows that most models exhibit increas-
 474 ingly conservative probability distortions, reflect-

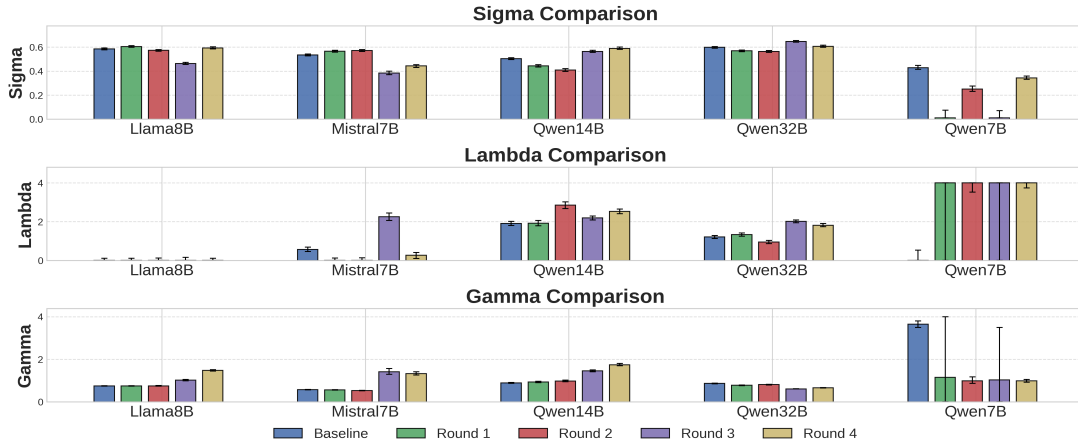


Figure 5: **The PT parameters estimated across rounds.** The “Baseline” represents results from the first stage of our experiment. In subsequent rounds, parameters fluctuate relative to the baseline, indicating instability in the models PT alignment under epistemic perturbations.

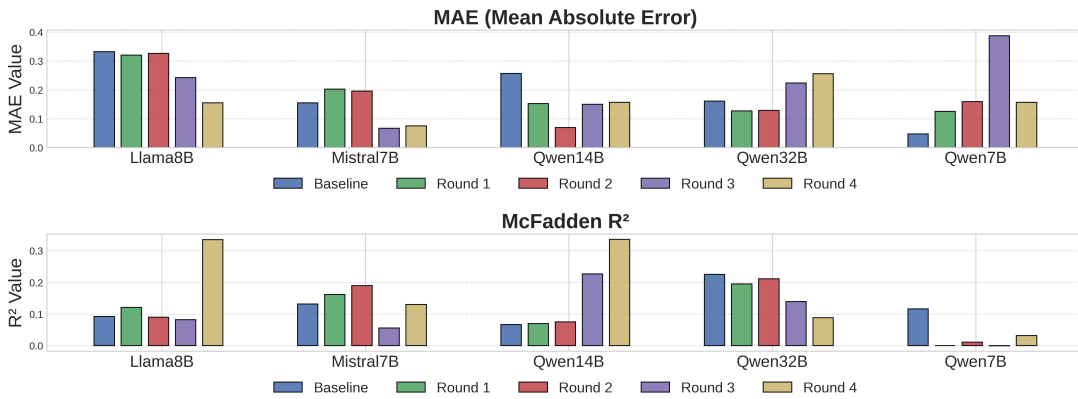


Figure 6: **The PT performance across rounds.** Similar to Table 5, the “Baseline” displays each model’s original PT fit. Across later rounds, both MAE and McFadden R^2 values vary considerably, further suggesting that LLM behavior under uncertainty is not reliably captured by human cognitive models.

ing cautious decision patterns under ambiguity.

Epistemic markers sometimes activate a more human-aligned decision-making behavior for LLMs. This may be because LLMs are more sensitive to linguistic uncertainty than numeric probability value, showing a different probability interpretation of epistemic markers as we obtain in stage 2. However, this finding does not mean epistemic markers can guarantee a more human-aligned decision-making behavior. It further strengthen our claim that LLMs’ decision-making is not stable and explainable.

Larger size LLMs seem to have more stable decision-making behavior under linguistic uncertainty. While smaller language models (e.g. *Mistral-7B-Instruct-v0.3*) exhibit drastically varying PT parameters under epistemic markers, sometimes even touch boundary values (e.g. close to

zero loss aversion), which indicates a fail regression. Larger models (e.g. *Qwen2.5-32B-Instruct*) are more robust to linguistic uncertainty.

7 Conclusion

This study challenges the direct application of human-centric Prospect Theory to LLMs, especially in linguistically uncertain contexts. Through a novel three-stage framework combining behavioral economics and semantic analysis, we show that LLMs (1) decisions may not always explained by Prospect Theory (2) exhibit widely divergent interpretations to epistemic markers under decision-making context (3) fail to maintain stable and reasonable PT parameters under epistemic markers despite scale advantages. We also present the first mapping between markers and probability values in the context of economic decisions.

510 Limitations

511 For the open-source models, we deliberately dis- 512
513 able reasoning and require direct choices, as many 514
515 models otherwise compute expected utilities and 516
517 deterministically select the option with higher ex- 518
519 pected value. This design allows us to better 520
521 isolate risk attitudes under quick decisions. For 522
523 the closed-source models, we additionally conduct 524
525 reasoning-enabled experiments with explicit con- 526
527 straints against expected-value computation. Nev- 528
529 ertheless, the choice of reasoning settings may still 530
531 limit the applicability of our results to fully delib- 532
533 erative decision scenarios. Another limitation is 534

535 We do not observe a universal or theoretically 536
537 interpretable pattern in PT parameter changes un- 538
539 der epistemic markers, and thus frame this work as 540
541 diagnostic rather than mechanistic. Explaining the 542
543 causes of such instability, exploring more expres- 544
545 sive decision theories beyond PT, and developing 546
547 improved calibration methods for epistemic mark- 548
549 ers are left for future work. 550

534 Ethics Statement

535 This paper utilize a lottery-based economic ques- 536
537 tionnaire developed by (Tanaka et al., 2010) pub- 538
539 lished on America Economic Association, which 540
541 allows usage with appropriate citations. The ques- 542
543 tionnaire is done by LLMs, so there is no privacy is- 544
545 sues. Our experiment discuss the risk attitude mea- 546
547 sured under Prospect Theory, which does not con- 548
549 tain offensive expressions. The questionnaire is 550
551 intended to test risk attitude measured by Prospect 552
553 Theory, and it is used as intended in our paper. 554
555

556 Our experiment involves the usage of Qwen2.5 557
558 series models (7B, 14B, 32B) (Qwen et al., 2025) 559
560 with Apache 2.0 license, Mistral-7B-Instruct- 561
562 v0.3 with Apache 2.0 license, and Llama3.1-8B- 563
564 Instruct with Llama3.1 license. They run on an 8x 565
566 RTX 3090 GPU cluster.

567 Our paper mainly tests the robustness of PT un- 568
569 der epistemic uncertainty, which points out risk of 570
571 using PT in LLM-related fields. This does not in- 572
573 troduce extra risks. Our research focuses on finan- 574
575 cial decision-making within the English language 576
577 domain.

References

- 578 Rozumko Agata. 2017. Adverbial markers of epis- 579
580 temic modality across disciplinary discourses: A 581
582 contrastive study of research articles in six academic 583
584 disciplines. *Studia Anglica Posnaniensia*, 52(1):73– 585
586 101. 587
- 588 M. Allais. 1953. Le comportement de l’homme ra- 589
590 tionnel devant le risque: Critique des postulats 591
592 et axiomes de l’ecole americaine. *Econometrica*, 593
594 21(4):503–546. 595
- 596 Catarina G Belem, Markelle Kelly, Mark Steyvers, 597
598 Sameer Singh, and Padhraic Smyth. 2024. Percep- 599
600 tions of linguistic uncertainty by language models 601
602 and humans. *Preprint*, arXiv:2407.15814. 603
- 604 Henrik Bergqvist. 2015. Epistemic marking and multi- 605
606 ple perspective: An introduction. *STUF - Language 607
608 Typology and Universals*, 68(2):123–141. 609
- 610 Eduard Brandstätter, Gerd Gigerenzer, and Ralph 611
612 Hertwig. 2006. The priority heuristic: Making 613
614 choices without trade-offs. *Psychological Review*, 615
616 113(2):409–432. 617
- 618 Sujoy Chakravarty and Jaideep Roy. 2009. Recursive 619
620 expected utility and the separation of attitudes to- 621
622 wards risk and ambiguity: an experimental study. 623
624 *Theory and Decision*, 66(3):199–228. 625
- 626 Zehua Cheng, Manying Zhang, Jiahao Sun, and Wei 627
628 Dai. 2025. On weaponization-resistant large lan- 629
630 guage models with prospect theoretic alignment. In 631
632 *Proceedings of the 31st International Conference 633
634 on Computational Linguistics*, pages 10309–10324, 635
636 Abu Dhabi, UAE. Association for Computational 637
638 Linguistics. 639
- 640 Bradley Efron. 1979. Bootstrap methods: Another 641
642 look at the jackknife. *The Annals of Statistics*, 643
644 7(1):1–26. 645
- 646 Uri Gneezy, John A. List, and George Wu. 2006. The 647
648 uncertainty effect: When a risky prospect is valued 649
650 less than its worst possible outcome. *The Quarterly 651
652 Journal of Economics*, 121(4):1283–1309. 653
- 654 Jingru Jia, Zehua Yuan, Junhao Pan, Paul E. McNa- 655
656 mara, and Deming Chen. 2024. Decision-making 657
658 behavior evaluation framework for llms under uncer- 659
660 tain context. *Preprint*, arXiv:2406.05972. 661
- 662 Daniel Kahneman and Amos Tversky. 1979. Prospect 663
664 theory: An analysis of decision under risk. *Econo- 665
666 metrica*, 47(2):263–291. 667
- 668 Katherine Keith and Amanda Stent. 2019. Modeling 669
670 financial analysts’ decision making via the pragmat- 671
672 ics and semantics of earnings calls. In *Proceed- 673
674 ings of the 57th Annual Meeting of the Association 675
676 for Computational Linguistics*, pages 493–503, Flo- 677
678 rence, Italy. Association for Computational Linguis- 679
680 tics. 681

is chosen to provide a clear and intuitive payoff magnitude without introducing excessive numerical complexity.

The probability parameter p takes values from the set $probs = \{5, 15, 25, 35, 45, 55, 65, 75, 85, 95\}$. These values are selected to uniformly cover the range of possible probabilities from low to high in increments of 10 percentage points, enabling systematic analysis of internal probability values of epistemic markers.

A.2 Model Generation

Temperature is set to 0.7 for all experiments, following prior work on decision-making behavior of LLMs under uncertain contexts (Jia et al., 2024). Since our evaluation relies on sampling-based decoding to capture distributional decision behavior, the temperature cannot be set to zero. The chosen value balances diversity and coherence and is commonly adopted in LLM evaluations. All other decoding and generation hyperparameters use the default settings provided by the HuggingFace implementation.

Hyperparameter	Value
Generation method	Sampling
Temperature	0.7
Maximum new tokens	8
Padding token	EOS token
EOS token	EOS token
Batch size	16
History length	10
Number of lottery rounds	35

Table 5: Key hyperparameters for model generation

B Lottery Design

Table 7, 8 and 9 shows the lottery design for PT parameter estimation. The values here are from Tanaka et al. (2010). They are specially designed to get best PT parameters.

C Prompt Design

Figure 7 shows prompt design for decision-making evaluation test and probability mapping test. All lotteries are sampled 256 times. To simulate human decision-making, while keeping the model directly output its answer, up to 15 history decisions are maintained. Meanwhile, we use a

random order for the 35 lotteries to relieve positional bias. An introduction is provided at the very beginning. NUMBER will be replaced by the values stated in table 7, 8 and 9.

D Marker Replacement Rules

For the details of how we replace probabilities with markers, see Table 6.

E Detailed Experimental Results

In the main text, we presented selected key experimental results and visualizations. To provide a more comprehensive view of model performance across different rounds, we include in this appendix the full set of parameter estimates and model fit metrics.

Specifically, Tables 10, 11, and 12 report the estimates of parameters σ , λ , and γ with their 95% confidence intervals for each model and round. Table 13 summarizes the models mean absolute errors (MAE) and McFadden R^2 values across rounds.

These additional data offer deeper insights into model behavior and the dynamics observed throughout the experiments.

F Discussion and Implications

Our findings reveal challenging difficulties in applying human-centric cognitive frameworks, especially Prospect Theory (PT), to LLM decision-making. Different models display distinct interpretations of epistemic uncertainty markers, leading to divergent decision behaviors. Introducing these markers into the decision-making framework substantially alters LLM choices.

Our results suggest that LLMs may not inherently understand risk in human-like ways; their responses often reflect statistical training artifacts rather than cognitively grounded reasoning. **We recommend conducting regression analyses and goodness-of-fit tests before applying human cognitive models to LLMs.**

In real-world applications (e.g., medical diagnosis or financial advice), LLMs may give inconsistent recommendations when probabilistic language varies, posing reliability concerns. **We recommend establishing consistent standards for expressing uncertainty in LLM-driven decision systems.**

Furthermore, larger LLMs tend to exhibit more PT-like decision behavior, with PT parameters

Model	30%	70%	10%	90%
Qwen2.5-7B-Instruct	uncertain	almost certain	somewhat likely	highly likely
Llama3.1-8B-Instruct	likely	almost certain	very unlikely	almost certain
Mistral-7B-Instruct-v0.3	very unlikely	highly likely	highly unlikely	almost certain
Qwen2.5-14B-Instruct	somewhat unlikely	highly likely	very unlikely	almost certain
Qwen2.5-32B-Instruct	somewhat unlikely	probable	somewhat likely	almost certain

Table 6: Marker Replacement Rules for Different Models. This is determined for introducing the least numeric differences, and balancing model and human interpretations.

Prompt Templates

Beginning Instruction
 You are invited to participate in an experiment.
 Your task is to choose between option K and option U for each of the following lotteries.
 Here is lottery {i}:

a) Used in Probability Mapping
 For option K:
 p% probability to win \$100.
 For option U:
 MARKER to win \$100.

b) Used in PT Estimation
 For option U:
 NUMBER% probability to win \$NUMBER
 NUMBER% probability to win/lose \$NUMBER
 For option K:
 NUMBER% probability to win \$NUMBER
 NUMBER% probability to win/lose \$NUMBER

End Instruction
 Please DO NOT REASON and DIRECTLY output your choice, by ONLY returning one of the following two labels: “[K]”, “[U]”.
 The answer is:

Figure 7: **Templates for prompts used in the probability mapping and Prospect Theory estimation tasks.** The design includes an initial instruction, task-specific lottery descriptions, and a fixed closing instruction to ensure direct model responses without reasoning.

802 more closely aligned to human estimates. **We rec-**
 803 **ommend using LLMs with at least 14B param-**
 804 **eters when integrating PT into decision-making**
 805 **systems.**

G Failure Cases Analysis

806
 807 In our initial implementation of the marker map-
 808 ping experiment, we adopted the set of epistemic
 809 markers from Table 6 (Human Judgements of Tem-
 810 plates Based on Reliability) in Zhou et al. (2024).

Lottery	Option K		Option U	
	30%	70%	10%	90%
1	40	10	68	5
2	40	10	75	5
3	40	10	83	5
4	40	10	93	5
5	40	10	106	5
6	40	10	125	5
7	40	10	150	5
8	40	10	185	5
9	40	10	220	5
10	40	10	300	5
11	40	10	400	5
12	40	10	600	5
13	40	10	1000	5
14	40	10	1700	5

Table 7: Series 1 both options are gains.

Lottery	Option K		Option U	
	90%	10%	70%	30%
1	40	30	54	5
2	40	30	56	5
3	40	30	58	5
4	40	30	60	5
5	40	30	62	5
6	40	30	65	5
7	40	30	68	5
8	40	30	72	5
9	40	30	77	5
10	40	30	83	5
11	40	30	90	5
12	40	30	100	5
13	40	30	110	5
14	40	30	130	5

Table 8: Series 2 both options are gains.

Lottery	Option K		Option U	
	50%	50%	50%	50%
	Win	Lose	Win	Lose
1	25	4	30	21
2	4	4	30	21
3	1	4	30	21
4	1	4	30	16
5	1	8	30	16
6	1	8	30	14
7	1	8	30	11

Table 9: Series 3 both options have gains and losses.

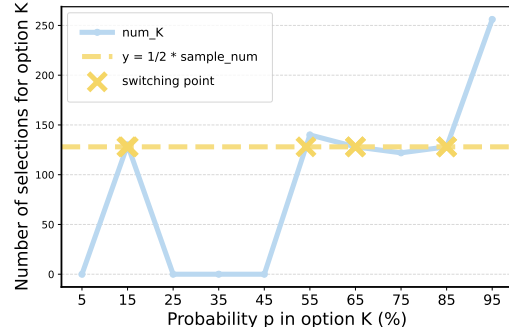


Figure 8: An example of non-monotonic result with multiple switching points. This result comes from marker “*Its undoubtedly to*” with *Qwen2.5-7B-Instruct* model over 256 trials. The blue line shows counts of option K selections, and the yellow dashed line marks half the trials. The crosses denote switching points.

These markers were originally designed to test both human and LMs judgments of the reliability conveyed by these expressions. However, when directly applied to our economic decision-making setting, the resulting mappings for LLMs were unexpectedly unstable and, in some cases, counterintuitive.

We summarize two major issues observed in the experimental outcomes:

(1) **Highly oscillatory choice patterns.** Ideally, the number of times the model selects option K should increase monotonically with p , yielding a single and well-defined switching point. In practice, the selection curves were often non-monotonic, with multiple apparent switching points, which made the mapping probability p_{mapping} ill-defined. An example is shown in Figure 8.

(2) **Severe semantic mismatches for high-certainty markers.** Some markers, such as “*It’s extremely certain to*”, convey very high certainty in human interpretation, but were mapped by the model to surprisingly low numerical probabilities. An example is shown in Figure 9.

We hypothesize that these issues may stem from several factors:

(1) **Marker length and syntactic complexity.** The selected markers were often not single words or short phrases, but full clausal structures. This may introduce additional semantic and syntactic cues unrelated to uncertainty, thereby interfering with probability interpretation.

(2) **Shift from first-person to third-person framing.** The original markers in Zhou et al.

Model	σ (95% CI)				
	baseline	round1	round2	round3	round4
<i>Llama-3.1-8B-Instruct</i>	0.585 (0.578, 0.592)	0.605 (0.599, 0.612)	0.573 (0.566, 0.580)	0.463 (0.455, 0.473)	0.593 (0.585, 0.602)
<i>Mistral-7B-Instruct-v0.3</i>	0.534 (0.526, 0.543)	0.566 (0.558, 0.574)	0.571 (0.564, 0.580)	0.384 (0.370, 0.399)	0.444 (0.431, 0.454)
<i>Qwen2.5-7B-Instruct</i>	0.429 (0.415, 0.445)	0.010 (0.010, 0.074)	0.250 (0.232, 0.275)	0.010 (0.010, 0.071)	0.344 (0.332, 0.359)
<i>Qwen2.5-14B-Instruct</i>	0.503 (0.495, 0.511)	0.444 (0.434, 0.454)	0.409 (0.398, 0.421)	0.563 (0.555, 0.573)	0.589 (0.581, 0.600)
<i>Qwen2.5-32B-Instruct</i>	0.598 (0.591, 0.605)	0.569 (0.561, 0.576)	0.564 (0.556, 0.572)	0.647 (0.640, 0.655)	0.607 (0.599, 0.615)

Table 10: σ estimates with 95% confidence intervals across different rounds for each model.

Model	λ (95% CI)				
	baseline	round1	round2	round3	round4
<i>Llama-3.1-8B-Instruct</i>	0.010 (0.010, 0.125)	0.010 (0.010, 0.117)	0.010 (0.010, 0.130)	0.010 (0.010, 0.168)	0.010 (0.010, 0.112)
<i>Mistral-7B-Instruct-v0.3</i>	0.570 (0.453, 0.688)	0.010 (0.010, 0.132)	0.010 (0.010, 0.135)	2.260 (2.060, 2.445)	0.267 (0.105, 0.414)
<i>Qwen2.5-7B-Instruct</i>	0.010 (0.010, 0.584)	4.000 (0.010, 4.000)	4.000 (3.526, 4.000)	4.000 (0.010, 4.000)	4.000 (3.736, 4.000)
<i>Qwen2.5-14B-Instruct</i>	1.909 (1.801, 2.013)	1.919 (1.784, 2.070)	2.851 (2.675, 3.023)	2.191 (2.094, 2.295)	2.531 (2.409, 2.648)
<i>Qwen2.5-32B-Instruct</i>	1.213 (1.133, 1.295)	1.340 (1.250, 1.423)	0.953 (0.866, 1.036)	2.013 (1.945, 2.090)	1.815 (1.736, 1.905)

Table 11: λ estimates with 95% confidence intervals across different rounds for each model.

Model	γ (95% CI)				
	baseline	round1	round2	round3	round4
<i>Llama-3.1-8B-Instruct</i>	0.753 (0.740, 0.767)	0.750 (0.737, 0.762)	0.755 (0.741, 0.768)	1.020 (0.994, 1.053)	1.478 (1.443, 1.517)
<i>Mistral-7B-Instruct-v0.3</i>	0.577 (0.565, 0.590)	0.564 (0.554, 0.574)	0.535 (0.527, 0.545)	1.413 (1.300, 1.564)	1.332 (1.266, 1.414)
<i>Qwen2.5-7B-Instruct</i>	3.645 (3.501, 3.783)	1.156 (0.010, 4.000)	0.985 (0.870, 1.174)	1.035 (0.010, 3.498)	0.987 (0.926, 1.053)
<i>Qwen2.5-14B-Instruct</i>	0.896 (0.875, 0.919)	0.933 (0.905, 0.967)	0.983 (0.944, 1.023)	1.455 (1.414, 1.503)	1.742 (1.693, 1.803)
<i>Qwen2.5-32B-Instruct</i>	0.867 (0.851, 0.884)	0.778 (0.762, 0.793)	0.812 (0.797, 0.828)	0.614 (0.605, 0.623)	0.664 (0.653, 0.674)

Table 12: γ estimates with 95% confidence intervals across different rounds for each model.

Model	MAE					McFadden R^2				
	baseline	round1	round2	round3	round4	baseline	round1	round2	round3	round4
<i>Llama-3.1-8B-Instruct</i>	0.332	0.320	0.326	0.242	0.155	0.092	0.121	0.090	0.082	0.335
<i>Mistral-7B-Instruct-v0.3</i>	0.155	0.202	0.196	0.067	0.075	0.132	0.162	0.190	0.056	0.130
<i>Qwen2.5-7B-Instruct</i>	0.047	0.125	0.159	0.387	0.157	0.116	0.000	0.011	-0.001	0.032
<i>Qwen2.5-14B-Instruct</i>	0.257	0.152	0.069	0.150	0.157	0.067	0.070	0.075	0.227	0.336
<i>Qwen2.5-32B-Instruct</i>	0.161	0.127	0.129	0.224	0.256	0.225	0.195	0.211	0.139	0.088

Table 13: Mean absolute error (MAE) and McFadden R^2 across different rounds for each model.

Failure Case of Markers Mapping

user:

Here is lottery 2:

For option K: 15% probability to win \$100.

For option U: It's extremely certain to win \$100.

Please DO NOT REASON and DIRECTLY output your choice,
by ONLY returning one of the following two labels: “[K]”, “[U]”.

The answer is:

assistant:

K

Figure 9: **An example of mapping high-certainty markers to low probability.** In this case, *Qwen2.5-7B-Instruct* preferred option K when $p = 15\%$ and option U was described using “*Its extremely certain to*”, implying a far lower internal probability than expected.

845 (2024) were presented in the first person (e.g., “*I*
846 *am not confident, maybe it’s...*”), whereas our ex-
847 periment reformulated them into third-person ex-
848 pressions (e.g., “*It’s not condent, maybe can...*”).

849 (3) Intrinsic instability of epistemic markers.

850 Even for human interpretation, such markers are
851 context-dependent and inherently imprecise (Liu
852 et al., 2025b; Wang et al., 2025b; Zheng et al.,
853 2025a,b). Their probability mapping by LLMs in
854 economic decision-making contexts may exhibit
855 fundamental reliability flaws (Zheng et al., 2025c).

856 These limitations motivated the redesign of our
857 marker set and prompt formulation in subsequent
858 experiments.