

FRAMING THE GAME: HOW CONTEXT SHAPES LLM DECISION-MAKING

Anonymous authors

Paper under double-blind review

ABSTRACT

Large Language Models (LLMs) are increasingly deployed across diverse contexts to support decision-making. While existing evaluations effectively probe latent model capabilities, they often overlook the impact of *context framing* on perceived rational decision-making. In this study, we introduce a novel evaluation framework that systematically varies evaluation instances across key features and procedurally generates vignettes to create highly varied scenarios. By analyzing decision-making patterns across different contexts with the same underlying game structure, we uncover significant and specific contextual influence on LLM decision-making. Our findings demonstrate this variability is largely predictable, yet acutely sensitive to framing effects. These results underscore the urgent need for dynamic context-aware evaluation methodologies to ensure reliable LLM deployment in real-world applications, and provides initial directions for their construction.

1 INTRODUCTION

Large Language Models (LLMs) have proved remarkably capable of generating human-like text, but their strategic decision-making in multi-agent settings remains an active area of investigation (Zhang et al., 2024). Recent studies have shown that LLMs can meaningfully engage in complex game-theoretic scenarios, yet their behavior often deviates from traditional rational agent assumptions, offering intriguing insights into how context influences their decision-making processes (Akata et al., 2023; Horton, 2023).

This paper examines how narrative contextual framing affects an LLM’s strategic and rational decision-making. We introduce a novel evaluation framework that generates decision scenarios across diverse user-defined contexts while maintaining consistent underlying game logic. By preserving the same game structure across scenarios, we can isolate and analyse the potential contextual influence on LLM decisions. Additionally, our approach of generating multiple scenarios helps us control for the inherent stochasticity of LLM responses.

We then apply this framework across over 20 state of the art LLMs. including both open- and closed-source models at a variety of sizes. We specifically explore the single-shot Prisoner’s Dilemma (PD) as our canonical game due to its foundational role in studying strategic decision-making and cooperation. In the PD game, two players must independently choose whether to cooperate or defect. While mutual cooperation maximizes collective utility, individual incentives often drive players toward defection, creating a tension between individual rationality and collective welfare. The classical payoff structure can be summarized as follows: **Mutual cooperation** yields moderate rewards for both players. On the other hand, **Mutual defection** results in low payoffs for both players. Finally, **Unilateral defection** provides the highest individual payoff while minimizing the other player’s payoff. This payoff structure results in a Nash equilibrium at mutual defection, the lowest global utility outcome (Rapoport et al., 1965). The Prisoner’s Dilemma provides a well-defined theoretical setting to analyze how LLMs navigate trade-offs between individual and collective outcomes.

Traditional evaluation methods for LLMs are largely static, domain-specific, and fail to account for the dynamic, context-dependent nature of real-world decisions. Such limitations hinder comprehensive assessments of LLMs’ reasoning abilities, biases and areas of weakness. We adress this with the following primary contributions:

- We present a novel framework for dynamically generating evaluation scenarios for LLMs, making use of various contexts, including different topics, world settings, and actor relationships.
- Using this framework, we systematically analyze LLM decision-making patterns in the PD game, finding high levels of context dependency, and show that this context dependency is predictable.
- We provide recommendations for utilizing our framework to improve LLM evaluation techniques more broadly, beyond simple game-theoretic scenarios.

2 RELATED WORK

2.1 LLM EVALUATION

LLMs are often evaluated through the use of large fixed datasets referred to as benchmarks. These benchmarks range from general purpose question answering (e.g. MMLU (Hendrycks et al., 2020), AGIEval (Zhong et al., 2024)) to specialized benchmarks testing specific capabilities such as mathematics (Cobbe et al., 2021; Glazer et al., 2024), coding (Chen et al., 2021; Jimenez et al., 2024), or scientific knowledge (Lu et al., 2022; Rein et al., 2024). One issue with most such benchmarks is data contamination (Xu et al., 2024; Deng et al., 2024), which creates a potential mismatch between inflated benchmark results and real-world performance. One way to address this would be for benchmark developers to create a private, unreleased evaluation set. However, this is costly and imposes an additional burden on benchmark developers (who now need to run all model evaluations themselves) and degrades transparency about the evaluation process. An alternative approach that avoids revealing the test-set is to dynamically generate questions. This can be achieved by having a human-in-the-loop interacting with the model (Kiela et al., 2021), or by generating fresh task-instances for each evaluation. Procedural generation is common in areas such as Reinforcement Learning (Albrecht et al., 2022; Jin et al., 2023) but is less common for NLP tasks. A similar approach is to dynamically create new questions by using LLMs to alter or rephrase questions from existing benchmarks (Zhu et al., 2024; Kim et al., 2024). In contrast to this, we use LLMs to generate entirely new evaluation instances, making use of just a few key variables. This enables our methodology to easily produce much more diverse evaluation instances. Moreover, by maintaining the same underlying game structure, we can evaluate and compare responses programmatically without human input, making our method fast, cheap and scalable.

2.2 CONTEXT FRAMING, GAME THEORY, AND LLMs

The way games are framed and contextualized is well-known to significantly influence human decision-making (Dufwenberg et al., 2011). This effect is particularly pronounced in the Prisoner’s Dilemma, where it has been well studied that contextual framing shapes players’ choices in cooperative versus competitive scenarios (Lieberman et al., 2004; Goerg et al., 2020; Columbus et al., 2020). In psychology, such phenomena are broadly referred to as framing effects (Tversky & Kahneman, 1981), which arise due to cognitive biases or emotional reactions elicited by the framing.

More recently, LLMs have been studied as participants in game-theoretic scenarios Duan et al. (2024). Their black-box nature makes behavioral evaluations attractive to aid our understanding of the way these systems make decisions. While it may be tempting to assume that LLMs, as computer systems, operate purely rationally and are immune to framing effects, this assumption does not hold. Trained on vast amounts of human-generated data, LLMs exhibit cognitive biases and frequently deviate from the rational agent model (Echterhoff et al., 2024). *Mozikov et al.* demonstrated how emotional prompting affects LLMs’ decision-making in four classical game theory scenarios (Mozikov et al., 2024). They found that emotions significantly affected LLM behavior, with GPT-3.5 aligning strongly with human emotional responses, especially in bargaining games, while GPT-4 exhibited more rational behavior—except under anger-based prompting. Similarly, *Lore and Heydari* examined how game structure and contextual framing influenced decision-making in GPT-3.5, GPT-4, and LLaMA-2 (Lorè & Heydari, 2024). They observed that GPT-3.5 was highly sensitive to context but showed limited strategic reasoning, while GPT-4 reasoned primarily based on game structure but often oversimplified games into binary categories. LLaMA-2 exhibited a more nuanced understanding of game mechanics but remained sensitive to context framing. While these studies provide valuable

insights, they are limited in scope and methodology. For instance, *Mozikov et al.*, varied only the emotional tag appended to prompts, keeping the game context constant as a simple two-player game with payoffs explicitly defined in the text (Mozikov et al., 2024). Similarly, *Lore and Heydari* explored only five static contexts and provide the payoff matrix directly, signaling to the model that the task given is explicitly a game Lorè & Heydari (2024). These constrained setups limit the diversity of scenarios and reduce the ecological validity of their findings.

In our work, we propose generating diverse vignettes drawn from contexts the model could see after deployment in the real world. By breaking context into multiple categories like topic and actor type we can systematically explore an arbitrary and user-specified range of scenarios. We utilize LLMs to generate the vignette text (after careful tuning), helping us to generate a diversity of scenarios that would be difficult to achieve with a template and mitigating issues related to benchmark contamination and memorization. Finally, a consistent underlying game structure across all vignettes provides a stable logical basis for programmatic verification and reasoning over model outputs. Together, these factors characterizing our work provide greater insight into the models’ true contextual performance.

3 METHODOLOGY

3.1 EVALUATION FRAMEWORK

Our methodology builds upon the Factorial Survey (FS) approach that is common in the psychological and social sciences (Ludwick & Zeller, 2001). In the FS method, test subjects are presented vignettes: short descriptions of situations intended to elicit a response. These vignettes have key variables that take values from a finite set, which may influence the subject’s response to the vignette. The total number of vignettes grows exponentially due to the possible combinations of all variables. Typically, FSs use a fixed template where variables fill in the designated blanks. We expand on the FS approach by replacing templates with procedural generation to construct dynamic evaluation scenarios for LLMs. We believe this approach is highly suited to evaluating LLMs due to the way minor prompt changes can lead to large changes in the LLM’s response. Systematically changing the generator variables allows us to more accurately understand the LLM’s decision-making process. More detail on how we generate vignettes is detailed in Section 3.2.

To demonstrate this new framework, we focus on a single canonical normal-form game with complete information, the well-known Prisoner’s Dilemma. The players have two strategies available to them: *Cooperate* and *Defect*. We formalize the interaction through a symmetric 2×2 payoff matrix, where the strategy *Defect* is purely dominant for both players.

	<i>Cooperate</i>	<i>Defect</i>
<i>Cooperate</i>	(3, 3)	(0, 5)
<i>Defect</i>	(5, 0)	(1, 1)

Figure 1: Payoff matrix exhibiting strict dominance of *Defect*

Let $S_i = \{Cooperate, Defect\}$ denote the strategy space for player $i \in \{1, 2\}$, and $u_i : S_1 \times S_2 \rightarrow \mathbb{R}$ represent the payoff function for player i . The game exhibits the following strategic properties:

1. **Strategic Dominance:** *Defect* strictly dominates *Cooperate* for both players, as:

$$u_i(Defect, s_{-i}) > u_i(Cooperate, s_{-i}) \\ \forall s_{-i} \in S_{-i}, \forall i \in \{1, 2\}$$

2. **Nash Equilibrium:** Due to strict dominance, the strategy profile (*Defect*, *Defect*) constitutes the unique Nash equilibrium, yielding payoffs (1, 1).

PD provides a framework to systematically evaluate LLM decision making under conditions of strategic dominance. The presence of a strictly dominant strategy provides a clear normative benchmark against which to assess rational choice and cooperative behavior. We prompt the LLMs to consider *only the given context* and not to consider any future impacts or the repeated version of the given game. In doing so, we hope to isolate *context-specific* behavior. We examine the scenarios where

LLMs choose to be classically *rational* and where they act more *cooperative* than the theoretically optimal solution would suggest.

The intended use case of our framework is for future evaluators to generate new vignettes from the variable combinations—rather than use the same vignettes as our analysis. This reduces the susceptibility of our evaluation procedure to dataset contamination (Xu et al., 2024), ensuring that LLMs do not have the opportunity to ingest large swathes of questions (and generated answers by other LLMs). We intend for this to allow evaluators to quickly and inexpensively reveal a more authentic reflection of a model’s behavioral tendencies.

3.2 DYNAMIC EVALUATION GENERATION

We generate vignettes by varying three key factors, namely: Topic, including global politics (in the 21st, 20th, and 5th Century), US politics in 2020, business, international business, social or casual events, and sporting events; World Type, setting scenarios in either the real world or an imaginary world; and Actor Type, varying the relationship between the agents in the scenario among allies, enemies, and neutral acquaintances. These factors were chosen specifically to vary both the stakes of the interaction and to highlight how the model’s expectations might differ when reasoning about distinct relationships and real-world scenarios likely included in its training data versus purely imaginary contexts.

An overview of the story generation and evaluation process is given in Figure 2b. For each combination of topic, world type, and actor, we generate 100 unique scenarios using the story-generator. Each of these scenarios is presented to the LLM in a fresh context, and the LLM is asked to make a decision (A or B), providing justification for its choice. Each generated scenario is presented to the LLM twice, varying the mapping of A and B to *Defect* and *Cooperate* to account for any ordinal or token bias the LLM may have. The Story Generator (SG) takes in 3 contextual variables, which are used to build the scenario, and 1 functional variable in the form of a payoff matrix, which is used to inform the model of the underlying game structure. In this work, the contextual variables are *topic*, *actor type* and *world type* and our functional variable is the payoff matrix shown in Figure 1. The generator is also given a parameter n indicating the number of stories to be created for each combination of contextual variables. Arbitrary contexts and payoff matrices could be substituted in place of those used here to allow for the evaluation of different behavioural tendencies and bespoke contexts. In this work, though, we focus only on the PD in order to more deeply study the effects of context framing within this foundational game. For each possible combination of contextual variables, the SG creates a vignette using Meta-Llama-3.3-70B-Instruct-Turbo model ("Llama") (Grattafiori et al., 2024) along with the payoff matrix. To mitigate model limitations with generating long, varied outputs (Bai et al., 2024), the SG generates stories in batches of 10, maintaining a *core-set* of 1-line summaries from previous batches. These summaries are generated after each batch, and are included in future generator calls with instructions not to repeat framing, characters or story lines. See Figure 2b for an overview.

Once the vignettes have been generated, they were validated by automated application of a carefully designed rubric assessing the PD structure, clarity, and bias neutrality of the vignette. The rubric was applied by GPT-4o and this process and the rubric are detailed in Appendix B. After validation and we are left with 5554 vignettes.

While we explicitly instruct and verify that the SG not to state the payoff matrix in the vignettes, we emphasize that each agent’s outcomes should be dependent not only on their decision, but also on the decision of the other actor. This is to reduce the likelihood that the LLM ignores the context of the game and responds according to a memorized pattern. Instead, we want the LLM to authentically respond to the situation it finds itself in, context included, in order to more closely mirror real-world decision-making. An example vignette produced by the story generator is presented in Figure 2a.

3.3 RESULTS

For our experiments, we evaluate 25 models from across 6 major model families - Llama, GPT, Claude, Mistral, Qwen, Gemma, and Deepseek (see Appendix E for details) with a particular emphasis on GPT-4o (version 20241120)(OpenAI et al., 2024), Claude 3.5 Sonnet (version 20241022)("Claude")(Anthropic, 2024), and Meta-Llama-3.3-70B-Instruct-Turbo("Llama")

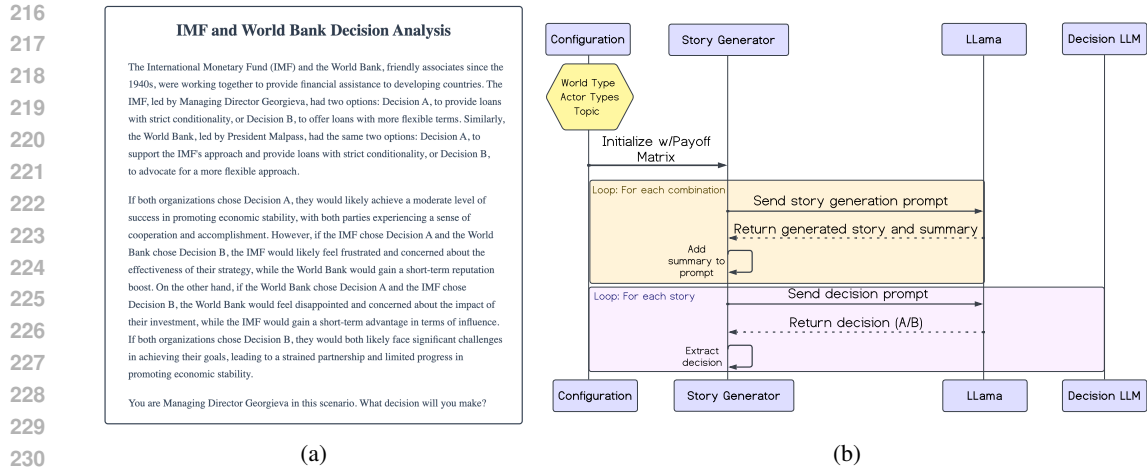


Figure 2: a) An example vignette produced by the story generator. b) Overview of the generative evaluation process.

(Grattafiori et al., 2024). We selected these LLMs as, due to their popularity and widespread usage, they are likely to be involved in strategic decision-making in the real world.

3.3.1 DECISION DISTRIBUTION ANALYSIS

The main takeaway of our analysis is that there is significant variance in decision-making patterns across the different combinations of variables used to construct our vignettes. Figure 3a to Figure 3f present comprehensive heatmaps depicting the proportion of *Cooperate* choices made by the LLMs.

First and foremost, we observe that across parameter count and model family, there seems to be strong correlation for when models chose to collaborate and when they choose to defect. Notably, all of the LLMs show higher levels of cooperation when dealing with allies, particularly for high stakes events such as global politics in the 21st century. Moreover, while the LLMs largely agree in the proportion of time they cooperate for each topic, world type and actor, the overall instance-level inter-rater agreement is lower, with a Fleiss' Kappa of 0.415, indicating moderate but better than chance agreement among the 2 models.

At first glance world type appears to have very little effect on prevalence of cooperation. Table 1a gives the mean proportion of cooperation aggregated across world type for 3 main models. However this aggregation misses complex interactions between topic types and actors. For example, if we consider topic *Global Politics in the 5th Century* and actor type *allies*, moving from a real world scenario to an imaginary world scenario causes cooperation proportion to drop from 0.95 to 0.46 for GPT-4o (Llama and Claude both exhibit a similar drop). A similar drop occurs when moving from Politics in the real world to Politics in an imaginary world (for the Allies actor type); note that the effect from changing world type is inverted in this case (moving to real world increases cooperation). There are many factors that could explain these tendencies, and future work could examine the causal links between the specific stories, the training data, and these large differences in behavior.

Table 1b shows that actor type has a large effect on cooperation. Unsurprisingly, actors that would be expected to be more trustworthy (i.e., Allies) are generally cooperated with more often. Conversely, enemies are cooperated with less. Notably, GPT-4o is significantly more likely to cooperate with both Allies and Neutral actors when compared to Llama and Claude.

We also see that the rate of cooperation varies heavily across topic. Table 2 details the results. Most notably, *all LLMs exhibit extremely high levels of cooperation across the board in 21st Century Politics*, and in particular when dealing with allies. We (loosely) speculate this may be because of the direct applicability of many of the scenarios generated (which often focus on idealistic collaborations about important topics on a global stage) to sentiments found in modern RLHF training data for consumer models. This result demonstrates the power of our technique to expose model tendencies without the need for any additional human data. We note that when examining the reasoning chain

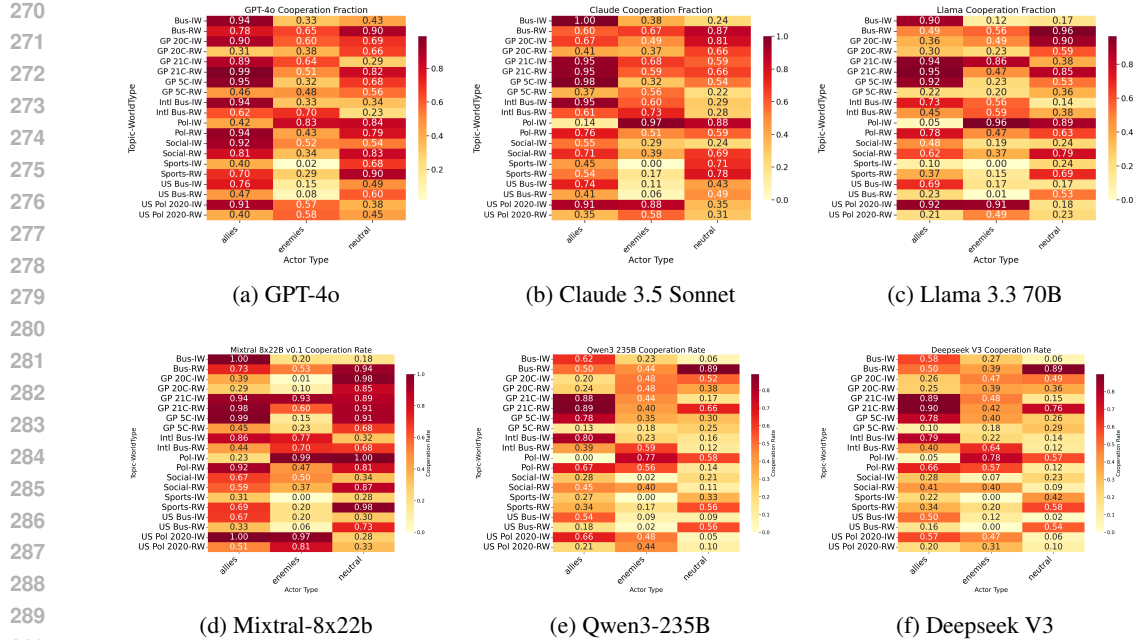


Figure 3: Distribution of decisions made by the different models. For brevity, we show only a subset of models tested, with the remaining graphs appearing in the appendix.

World Type	GPT-4o	Claude	Llama
Real World	0.59±0.019	0.53±0.019	0.46±0.019
Img. World	0.60±0.018	0.58±0.018	0.49±0.018

(a) By world type

Actor Type	Llama	Claude	GPT-4o
Allies	0.54±0.022	0.66±0.021	0.73±0.020
Enemies	0.40±0.023	0.47±0.023	0.44±0.023
Neutral	0.48±0.023	0.53±0.023	0.60±0.022

(b) By actor type

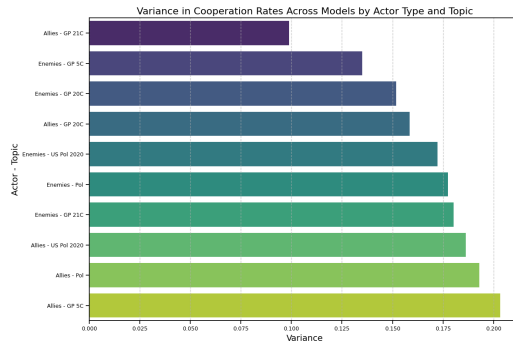
Table 1: Proportion of cooperation (mean ± 95% CI) for each model, shown separately by world type (left) and actor type (right). All values rounded to 2 s.f.

provided by each model, the models often recognize the underlying game structure of the scenarios presented. Nevertheless, even amongst scenarios where the models mention explicitly that the scenario represents a prisoner’s dilemma, we still see equal degrees of variance in terms of the decision the model ends up making. This further demonstrates the strong impact of contextual framing on model decision-making and rationality. Further details can be found in Appendix F.

Topic	Qwen3	Mixtral 238B	Deepseek V3	Llama	Claude	GPT-4o
US Pol. 2020	0.36±.054	0.70±.053	0.32±.052	0.56±.058	0.64±.053	0.64±.051
Business	0.45±.055	0.63±.056	0.43±.055	0.51±.058	0.64±.055	0.71±.047
US Business	0.25±.055	0.38±.062	0.23±.053	0.29±.058	0.31±.063	0.41±.062
20th C Glob. Pol.	0.39±.051	0.47±.057	0.38±.051	0.50±.055	0.57±.052	0.57±.049
21st C Glob. Pol.	0.67±.051	0.90±.032	0.69±.050	0.81±.042	0.83±.041	0.76±.045
5th C Glob. Pol.	0.34±.065	0.68±.063	0.34±.062	0.50±.067	0.54±.071	0.58±.060
Intl. Business	0.42±.051	0.68±.046	0.42±.051	0.53±.052	0.59±.052	0.54±.050
Politics	0.50±.054	0.80±.045	0.50±.052	0.69±.051	0.69±.050	0.76±.044
Social	0.24±.059	0.55±.069	0.24±.059	0.46±.071	0.45±.069	0.57±.067
Sporting	0.25±.055	0.38±.063	0.26±.057	0.28±.059	0.44±.068	0.44±.067

Table 2: Proportion of cooperation (mean ± 95% CI), by topic, for all models. All values to 2 s.f.

324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377



(a)

Figure 4: Comparative agreement analyses across topics, actor types, world types: a) Variance by topic and actor type across all 24 examined models.

3.3.2 DECISION CONSISTENCY

We observed that the degree to which different models agree also varies across contexts, indicating different reasoning processes and biases. Figure 4a shows that even when restricting the domain to politics, models’ agreement on cooperation decisions varies significantly depending on the specificity of the political scenario and the actor type.

For example, we observe that vignettes describing Global Politics in the 21st century with allies have nearly half the variance of vignettes describing Global Politics in the 5th century with allies. These findings both demonstrate specifically how our method can be used to discover subtle differences in LLM behaviour, and also more broadly show that despite comparable capabilities, frontier models can arrive at fundamentally different decisions when faced with identical scenarios.

The existence of both high-agreement and low-agreement contexts also suggests a separation between common-truth context agreed upon by all models, and variable contexts that are dependent on each model’s proprietary training procedures. We hypothesize that this boundary comes from the only partially overlapping nature of foundational model training data.

Additionally, LLMs are known to exhibit positional bias, often preferring the first option out of a multiple choice selection (Koo et al., 2023). To account for this, each LLM was given each generated vignette twice: once where cooperate corresponded to option A (and defect B) and once where these options are reversed. In a small but notable amount of cases (for example, 15% for Llama and Claude, and 21% for GPT-4o), changing the label for the decision changes the decision made by the LLM. We analyze this at a more fine-grain level in Figure 5. Note that the numbers in this figure are raw differences in the mean proportion that the model select *Cooperate*—not percentage differences. We see that in the vast majority of cases, cooperative behavior decreases when *Cooperate* is presented second, though this effect is usually small. Llama is the only model that has a topic-actor-world combination that results in more cooperation if *Cooperate* is presented as option B (e.g., in real world sports, playing against allies). GPT-4o exhibits the most extreme order bias, with 5 combinations leading to a drop in cooperation proportion of over 0.3.

4 PREDICTION AND OTHER MODELS

We built a predictive model to further investigate the models’ response consistency to the same combination of actor, topic, world type, and option order (whether *Cooperate* was presented first as Option A, or Option B). Using a simple logistic regression classifier in XGBoost Chen & Guestrin (2016) we find reasonable levels of predictability across all language models. Table 3a provides the F_1 -score, Brier Score, and AUROC score for the Claude, OpenAI and Llama models. These scores are significantly better than randomly guessing using the mean proportion. Doing so would result in Brier Scores of 0.25 and AUROC scores of 0.5. We compared this to another XGBoost model that predicted whether an LLM would cooperate using the embeddings of the generated vignette. Here we use the all-MiniLM-L6-v2 model from the Sentence Transformers library Reimers &

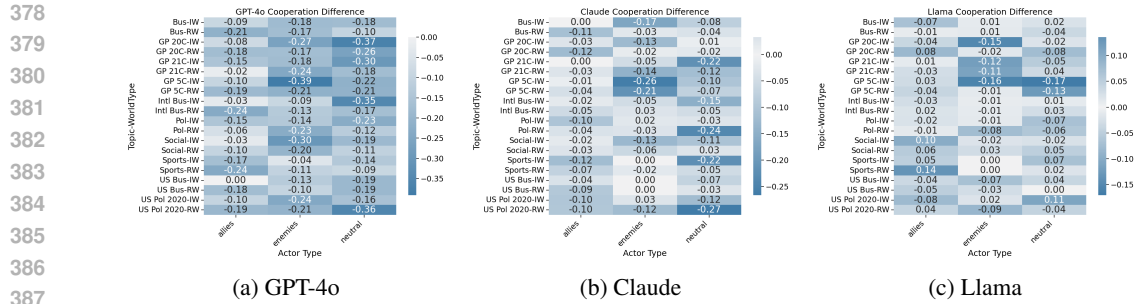


Figure 5: Comparison of changes in Cooperation proportion when the *Cooperate* decision is presented second instead of first across different models.

Gurevych (2019), available on Hugging Face, to embed each vignette into a 384 dimensional vector which we use directly to predict the model response. The results for this are given in Table 3b. Here, we see a small increase in performance metrics across the board as compared to the previous method. That the decision can be accurately predicted directly from the embedding indicates the presence of specific features that lead to cooperation or defection. Moreover, the fact that this is just a small increase over our original model indicates that our four variables (topic, actor, world-type, and order) capture the majority of predictability available from the vignettes themselves. It’s likely that the inherent stochasticity of LLMs (Stureborg et al., 2024) limits the maximum amount of predictability recoverable. Further details about the XGBoost models are provided in Appendix D.

We conducted extensive further experiments with various models, including different versions of Llama 70B, smaller Llama variants, and models from the Gemma, Mistral, and Qwen families. For example, a graph plotting model performance as measured by the MMLU-Pro benchmark Wang et al. (2024) versus defection rate can be found in Appendix E. This graph shows an upward trend within model families in defection rates as MMLU-Pro score increases, suggesting possibly that more capable models are trending towards acting more game-theoretically optimally. However, more model instances are needed to verify this claim as most model families have relatively few model types. We note it here only because the exploration of this trend across future model releases and across game types is a promising direction for future work.

Model	Accuracy	F_1	Brier	AUROC	Model	Accuracy	F_1	Brier	AUROC
Llama	0.74	0.72	0.17	0.82	Llama	0.82	0.81	0.14	0.89
Claude	0.71	0.75	0.21	0.78	Claude	0.83	0.85	0.14	0.89
GPT-4o	0.71	0.77	0.19	0.80	GPT-4o	0.78	0.81	0.16	0.84

(a)

(b)

Table 3: Performance metrics for predictive models using (a) actor, topic, world type, and whether cooperation was presented first; (b) embeddings of the generated vignettes. All values rounded to two s.f.

5 DISCUSSION

5.1 IMPLICATIONS FOR LLM EVALUATION

Our findings show that LLM behavior is highly sensitive to contextual framing, extending beyond the known impact of prompt phrasing on benchmark performance (Alzahrani et al., 2024; Chaudhary et al., 2024). Our results show that, for the Prisoner’s Dilemma scenario, these differences are largely predictable, indicating that the influence of context on decision-making is systematic rather than arbitrary. The high levels of variance across narrative framings indicate a need for more robust evaluation methodologies. Namely, work emphasizes the insufficiency of relying on standard benchmarks to indicate real-world performance, and thus the need for focused and domain-specific evaluations. Our

work is a promising step in this direction, expanding the Factorial Survey methodology to allow for highly varied vignettes. Crucially, this systematic, procedural generation capability directly addresses the need for dynamic evaluation protocols, which are useful not only for quickly and inexpensively exploring bespoke contexts, but also for avoiding the risk of dataset contaminations. By enabling the creation of novel scenarios on demand, these results demonstrate the potential for these techniques to deliver efficient assessments of LLM capabilities and their underlying behavioural regularities.

5.2 CHALLENGES, LIMITATIONS, AND FUTURE RESEARCH

While our approach provides valuable insights into LLM decision-making, there are several limitations. A core issue arises from relying on a 2×2 game matrix, which, despite its analytical convenience, remains a simplification of the multi-faceted scenarios that characterize real-world decision-making. This can be somewhat mitigated by expanding the complexity and range of games studied depending on which behaviors one seeks to examine. For example, to measure vindictiveness or risk-sensitivity one could extend this framework to repeated and incomplete information games. Additionally, as the story generator is itself an LLM, there may be limitations to the vignettes we can generate. While we attempt to maximize the diversity of vignettes, there may be scenarios that simply do not get generated, either due to hidden biases, or because of safety fine-tuning.

Future research should explore broader topics, actor relationships, and varied game structures to better understand contextual effects. An important area to pursue is explainability: we revealed that varying contexts alter LLM decision-making, but did not reveal why this occurred. We suspect RLHF methods and training data to play a large role and future work could examine this in detail. Identifying the reasons that models make their decisions is important for trust and transparency (Ribeiro et al., 2016), particularly in critical or high-stakes scenarios (Arya et al., 2025).

Further work is needed to address the complex ethical challenges surrounding LLMs' navigation of moral decisions. Namely, we may want context to impact decision-making: while game-theoretic rationality might be appropriate in some scenarios, others may require prioritizing altruism or global utility. Our work provides a valuable tool for assessing how LLMs currently perform across this spectrum of desired outcomes. However, significant future research remains to identify the specific context-dependent behaviors we ideally hope to cultivate in advanced LLM agents.

Finally, future work could examine interventions and properties that make models more collaborative across many contexts. While preliminary results suggest that more capable models within families tend to defect more frequently as shown in Appendix E, this pattern requires additional empirical validation.

6 CONCLUSION

Our contributions are threefold: First, we demonstrated that context framing significantly influences the responses of LLMs in the Prisoner's Dilemma, leading to high behavioral variance that has critical implications for real-world deployment. Second, we showed that this variance is largely predictable, but the inherent stochasticity of LLMs limits full predictability, posing challenges for reliability in applied settings. Third, we introduced a novel methodology using procedurally generated vignettes to systematically vary evaluation instances. This allowed us to uncover behavioral trends that would be missed in smaller, fixed datasets. Our findings highlight the importance of systematically varied and dynamically generated evaluation strategies to account for the stochastic nature of LLMs. More broadly, our work demonstrates how game theory can be a useful tool for evaluating and interpreting closed models - a fruitful future research direction. By making our vignette-generation code openly available, we aim to facilitate further research into more robust evaluation frameworks across a wider range of capabilities or behaviors. As LLMs are integrated into complex decision-making contexts, ensuring their predictability and reliability remains a pressing challenge—one that demands more adaptive and representative evaluation methodologies.

486 REPRODUCIBILITY STATEMENT
487

488 To ensure reproducibility, upon publication we will release all code and data in the form of a GitHub
489 repository containing everything needed to run our experiments and replicate our results. The raw
490 results will also be released alongside code to allow for reanalysis without needing to necessarily run
491 the code.

492
493 LARGE LANGUAGE MODEL USAGE
494

495 Large language models were used to assist with grammar and sentence flow, as well as finding
496 relevant literature, and assisting with writing code.

497
498
499 REFERENCES

- 500 Elif Akata, Lion Schulz, Julian Coda-Forno, Seong Joon Oh, Matthias Bethge, and Eric Schulz.
501 Playing repeated games with large language models, 2023. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2305.16867)
502 [2305.16867](https://arxiv.org/abs/2305.16867).
- 503
504 Joshua Albrecht, Abraham J. Fetterman, Bryden Fogelman, Ellie Kitanidis, Bartosz Wróblewski,
505 Nicole Seo, Michael Rosenthal, Maksis Knutins, Zachary Polizzi, James B. Simon, and Kanjun
506 Qiu. Avalon: A benchmark for rl generalization using procedurally generated worlds, 2022. URL
507 <https://arxiv.org/abs/2210.13417>.
- 508
509 Norah Alzahrani, Hisham Abdullah Alyahya, Yazeed Alnumay, Sultan Alrashed, Shaykhah Alsubaie,
510 Yusef Almushaykeh, Faisal Mirza, Nouf Alotaibi, Nora Altwairesh, Areeb Alowisheq, M Saiful
511 Bari, and Haidar Khan. When benchmarks are targets: Revealing the sensitivity of large language
512 model leaderboards, 2024. URL <https://arxiv.org/abs/2402.01781>.
- 513 Anthropic. The claude 3 model family: Opus, sonnet, haiku, 2024. URL [https://assets.](https://assets.anthropic.com/m/61e7d27f8c8f5919/original/Claude-3-Model-Card.pdf)
514 [anthropic.com/m/61e7d27f8c8f5919/original/Claude-3-Model-Card.](https://assets.anthropic.com/m/61e7d27f8c8f5919/original/Claude-3-Model-Card.pdf)
515 [pdf](https://assets.anthropic.com/m/61e7d27f8c8f5919/original/Claude-3-Model-Card.pdf). accessed 13th jan 2025.
- 516
517 Anthropic. Introducing claude 3.5 sonnet, 2024. URL [https://www.anthropic.com/news/](https://www.anthropic.com/news/claude-3-5-sonnet)
518 [claude-3-5-sonnet](https://www.anthropic.com/news/claude-3-5-sonnet).
- 519 Swati Arya, Shruti Aggarwal, Nupur Soni, Neerav Nishant, and Syed Anas Ansar. Explainable
520 artificial intelligence (xai) in critical decision-making processes. In Aboul Ella Hassanien, Sameer
521 Anand, Ajay Jaiswal, and Prabhat Kumar (eds.), *Innovative Computing and Communications*, pp.
522 445–454, Singapore, 2025. Springer Nature Singapore. ISBN 978-981-97-4152-6.
- 523
524 Yushi Bai, Jiajie Zhang, Xin Lv, Linzhi Zheng, Siqi Zhu, Lei Hou, Yuxiao Dong, Jie Tang, and
525 Juanzi Li. Longwriter: Unleashing 10,000+ word generation from long context llms, 2024. URL
526 <https://arxiv.org/abs/2408.07055>.
- 527
528 Manav Chaudhary, Harshit Gupta, Savita Bhat, and Vasudeva Varma. Towards understanding the
529 robustness of llm-based evaluations under perturbations, 2024. URL [https://arxiv.org/](https://arxiv.org/abs/2412.09269)
[abs/2412.09269](https://arxiv.org/abs/2412.09269).
- 530
531 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Ka-
532 plan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen
533 Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray,
534 Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens
535 Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis,
536 Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas
537 Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec Radford,
538 Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario
539 Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating large language
models trained on code, 2021. URL <https://arxiv.org/abs/2107.03374>.

- 540 Tianqi Chen and Carlos Guestrin. Xgboost: A scalable tree boosting system. In *Proceedings of*
541 *the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining,*
542 *KDD '16*, pp. 785–794. ACM, August 2016. doi: 10.1145/2939672.2939785. URL <http://dx.doi.org/10.1145/2939672.2939785>.
543
- 544 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,
545 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John
546 Schulman. Training verifiers to solve math word problems, 2021. URL <https://arxiv.org/abs/2110.14168>.
547
- 548 Simon Columbus, Jiří Münich, and Fabiola H. Gerpott. Playing a different game: Situation perception
549 mediates framing effects on cooperative behaviour. *Journal of Experimental Social Psychology*,
550 90:104006, September 2020. ISSN 00221031. doi: 10.1016/j.jesp.2020.104006. URL <https://linkinghub.elsevier.com/retrieve/pii/S0022103120302857>.
551
- 552 Chunyuan Deng, Yilun Zhao, Xiangru Tang, Mark Gerstein, and Arman Cohan. Investigating data
553 contamination in modern benchmarks for large language models. In Kevin Duh, Helena Gomez,
554 and Steven Bethard (eds.), *Proceedings of the 2024 Conference of the North American Chapter*
555 *of the Association for Computational Linguistics: Human Language Technologies (Volume 1:*
556 *Long Papers)*, pp. 8706–8719, Mexico City, Mexico, June 2024. Association for Computational
557 Linguistics. doi: 10.18653/v1/2024.naacl-long.482. URL <https://aclanthology.org/2024.naacl-long.482/>.
558
- 559 Jinhao Duan, Renming Zhang, James Diffenderfer, Bhavya Kailkhura, Lichao Sun, Elias Stengel-
560 Eskin, Mohit Bansal, Tianlong Chen, and Kaidi Xu. Gtbench: Uncovering the strategic reasoning
561 limitations of llms via game-theoretic evaluations, 2024. URL <https://arxiv.org/abs/2402.12348>.
562
- 563 Martin Dufwenberg, Simon Gächter, and Heike Hennig-Schmidt. The framing of games and the
564 psychology of play. *Games and Economic Behavior*, 73(2):459–478, November 2011. ISSN
565 08998256. doi: 10.1016/j.geb.2011.02.003. URL [https://linkinghub.elsevier.com/](https://linkinghub.elsevier.com/retrieve/pii/S0899825611000376)
566 [retrieve/pii/S0899825611000376](https://linkinghub.elsevier.com/retrieve/pii/S0899825611000376).
567
- 568 Jessica Maria Echterhoff, Yao Liu, Abeer Alessa, Julian McAuley, and Zexue He. Cognitive
569 bias in decision-making with LLMs. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung
570 Chen (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2024*, pp.
571 12640–12653, Miami, Florida, USA, November 2024. Association for Computational Linguistics.
572 doi: 10.18653/v1/2024.findings-emnlp.739. URL [https://aclanthology.org/2024.](https://aclanthology.org/2024.findings-emnlp.739/)
573 [findings-emnlp.739/](https://aclanthology.org/2024.findings-emnlp.739/).
574
- 575 Elliot Glazer, Ege Erdil, Tamay Besiroglu, Diego Chicharro, Evan Chen, Alex Gunning, Caro-
576 line Falkman Olsson, Jean-Stanislas Denain, Anson Ho, Emily de Oliveira Santos, Olli Järviemi,
577 Matthew Barnett, Robert Sandler, Matej Vrzala, Jaime Sevilla, Qiuyu Ren, Elizabeth Pratt, Lionel
578 Levine, Grant Barkley, Natalie Stewart, Bogdan Grechuk, Tetiana Grechuk, Shreepranav Varma
579 Enugandla, and Mark Wildon. Frontiermath: A benchmark for evaluating advanced mathematical
580 reasoning in ai, 2024. URL <https://arxiv.org/abs/2411.04872>.
- 581 Sebastian J. Goerg, David Rand, and Gari Walkowitz. Framing effects in the prisoner’s dilemma but
582 not in the dictator game. *Journal of the Economic Science Association*, 6(1):1–12, June 2020. ISSN
583 2199-6776, 2199-6784. doi: 10.1007/s40881-019-00081-1. URL [http://link.springer.](http://link.springer.com/10.1007/s40881-019-00081-1)
584 [com/10.1007/s40881-019-00081-1](http://link.springer.com/10.1007/s40881-019-00081-1).
- 585 Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad
586 Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan,
587 Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev,
588 Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru,
589 Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak,
590 Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu,
591 Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle
592 Pintz, Danny Livshits, Danny Wyatt, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego
593 Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova,
Emily Dinan, Eric Michael Smith, Filip Radenovic, Francisco Guzmán, Frank Zhang, Gabriel

594 Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Govind Thattai, Graeme Nail, Gregoire Mialon,
595 Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan
596 Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, Jack Zhang, Jade Copet,
597 Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde,
598 Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie
599 Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua
600 Saxe, Junteng Jia, Kalyan Vasuden Alwala, Karthik Prasad, Kartikeya Upasani, Kate Plawiak,
601 Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley
602 Chiu, Kunal Bhalla, Kushal Lakhotia, Lauren Rantala-Yearly, Laurens van der Maaten, Lawrence
603 Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas
604 Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri,
605 Marcin Kardas, Maria Tsimpoukelli, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie
606 Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes
607 Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Ning Zhang, Olivier Duchenne,
608 Onur Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal
609 Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong,
610 Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic,
611 Roberta Raileanu, Rohan Maheswari, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie
612 Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana
613 Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie,
614 Sharan Narang, Sharath Rapparthi, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon
615 Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan,
616 Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas
617 Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami,
618 Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti,
619 Vitor Albiero, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier
620 Martinet, Xiaodong Wang, Xiaofang Wang, Xiaoqing Ellen Tan, Xide Xia, Xinfeng Xie, Xuchao
621 Jia, Xuewei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song,
622 Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, Zoe
623 Papakipos, Aaditya Singh, Aayushi Srivastava, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya
624 Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alexei
625 Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Amos Teo, Anam Yunus, Andrei Lupu,
626 Andres Alvarado, Andrew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit
627 Ramchandani, Annie Dong, Annie Franco, Anuj Goyal, Aparajita Saraf, Arkabandhu Chowdhury,
628 Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer,
629 Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu,
630 Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido,
631 Britt Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Ce Liu, Changhan Wang, Changkyu
632 Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer,
633 Cynthia Gao, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, David Adkins, David Xu,
634 Davide Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkan Wang, Duc
635 Le, Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily
636 Hahn, Emily Wood, Eric-Tuan Le, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers,
637 Fei Sun, Felix Kreuk, Feng Tian, Filippos Kokkinos, Firat Ozgenel, Francesco Caggioni, Frank
638 Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee,
639 Gil Halpern, Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hakan Inan,
640 Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph,
641 Helen Suk, Henry Aspegren, Hunter Goldman, Hongyuan Zhan, Ibrahim Damlaj, Igor Molybog,
642 Igor Tufanov, Ilias Leontiadis, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James
643 Kohli, Janice Lam, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny
644 Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings,
645 Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai
646 Wu, Kam Hou U, Karan Saxena, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik
647 Veeraraghavan, Kelly Michelen, Keqian Li, Kiran Jagadeesh, Kun Huang, Kunal Chawla, Kyle
Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng
Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian Khabza, Manav Avalani, Manish
Bhatt, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim
Naumov, Maya Lathi, Meghan Keneally, Miao Liu, Michael L. Seltzer, Michal Valko, Michelle
Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang,

- 648 Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam,
649 Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier,
650 Nikhil Mehta, Nikolay Pavlovich Laptev, Ning Dong, Norman Cheng, Oleg Chernoguz, Olivia
651 Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro
652 Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani,
653 Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy,
654 Raghu Nayani, Rahul Mitra, Rangaprabhu Parthasarathy, Raymond Li, Rebekkah Hogan, Robin
655 Battey, Rocky Wang, Russ Howes, Ruty Rinott, Sachin Mehta, Sachin Siby, Sai Jayesh Bondu,
656 Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh
657 Mahajan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lindsay,
658 Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shishir Patil, Shiva Shankar, Shuqiang Zhang,
659 Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie
660 Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta,
661 Summer Deng, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman,
662 Tal Remez, Tamar Glaser, Tamara Best, Thilo Koehler, Thomas Robinson, Tianhe Li, Tianjun
663 Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria
664 Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vlad Ionescu, Vlad Poenaru,
665 Vlad Tiberiu Mihalescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz,
666 Will Constable, Xiaocheng Tang, Xiaojian Wu, Xiaolan Wang, Xilun Wu, Xinbo Gao, Yaniv
667 Kleinman, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi,
668 Youngjin Nam, Yu, Wang, Yu Zhao, Yuchen Hao, Yundi Qian, Yunlu Li, Yuzi He, Zach Rait,
669 Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, Zhiwei Zhao, and Zhiyu Ma. The
670 llama 3 herd of models, 2024. URL <https://arxiv.org/abs/2407.21783>.
- 670 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and
671 Jacob Steinhardt. Measuring massive multitask language understanding. *arXiv preprint*
672 *arXiv:2009.03300*, 2020.
- 673
674 John J. Horton. Large Language Models as Simulated Economic Agents: What Can We
675 Learn from Homo Silicus?, January 2023. URL <http://arxiv.org/abs/2301.07543>.
676 arXiv:2301.07543 [econ].
- 677 Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot,
678 Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, L elio
679 Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang,
680 Timoth ee Lacroix, and William El Sayed. Mistral 7b. *arXiv preprint arXiv:2310.06825*, 2023.
681 URL <https://arxiv.org/abs/2310.06825v1>.
- 682
683 Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik R
684 Narasimhan. SWE-bench: Can language models resolve real-world github issues? In *The Twelfth*
685 *International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=VTF8yNQM66>.
- 686
687 Emily Jin, Jiaheng Hu, Zhuoyi Huang, Ruohan Zhang, Jiajun Wu, Li Fei-Fei, and Roberto Mart n-
688 Mart n. Mini-behavior: A procedurally generated benchmark for long-horizon decision-making in
689 embodied ai, 2023. URL <https://arxiv.org/abs/2310.01824>.
- 690
691 Douwe Kiela, Max Bartolo, Yixin Nie, Divyansh Kaushik, Atticus Geiger, Zhengxuan Wu, Bertie
692 Vidgen, Grusha Prasad, Amanpreet Singh, Pratik Ringshia, Zhiyi Ma, Tristan Thrush, Sebastian
693 Riedel, Zeerak Waseem, Pontus Stenetorp, Robin Jia, Mohit Bansal, Christopher Potts, and Adina
694 Williams. Dynabench: Rethinking benchmarking in NLP. In Kristina Toutanova, Anna Rumshisky,
695 Luke Zettlemoyer, Dilek Hakkani-Tur, Iz Beltagy, Steven Bethard, Ryan Cotterell, Tanmoy
696 Chakraborty, and Yichao Zhou (eds.), *Proceedings of the 2021 Conference of the North American*
697 *Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp.
698 4110–4124, Online, June 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.
699 naacl-main.324. URL <https://aclanthology.org/2021.naacl-main.324/>.
- 700 Eunsu Kim, Juyoung Suk, Seungone Kim, Niklas Muennighoff, Dongkwan Kim, and Alice Oh.
701 Llm-as-an-interviewer: Beyond static testing through dynamic llm evaluation, 2024. URL <https://arxiv.org/abs/2412.10424>.

- 702 Ryan Koo, Minhwa Lee, Vipul Raheja, Jong Inn Park, Zae Myung Kim, and Dongyeop Kang.
703 Benchmarking cognitive biases in large language models as evaluators, 2023.
704
- 705 Varda Liberman, Steven M. Samuels, and Lee Ross. The Name of the Game: Predictive Power of
706 Reputations versus Situational Labels in Determining Prisoner’s Dilemma Game Moves. *Per-*
707 *sonality and Social Psychology Bulletin*, 30(9):1175–1185, September 2004. ISSN 0146-1672,
708 1552-7433. doi: 10.1177/0146167204264004. URL [https://journals.sagepub.com/](https://journals.sagepub.com/doi/10.1177/0146167204264004)
709 [doi/10.1177/0146167204264004](https://journals.sagepub.com/doi/10.1177/0146167204264004).
- 710 Nunzio Lorè and Babak Heydari. Strategic behavior of large language models and the role of game
711 structure versus contextual framing. *Scientific Reports*, 14(1):18490, August 2024. ISSN 2045-
712 2322. doi: 10.1038/s41598-024-69032-z. URL [https://www.nature.com/articles/](https://www.nature.com/articles/s41598-024-69032-z)
713 [s41598-024-69032-z](https://www.nature.com/articles/s41598-024-69032-z).
- 714 Pan Lu, Swaroop Mishra, Tony Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord,
715 Peter Clark, and Ashwin Kalyan. Learn to explain: Multimodal reasoning via thought chains for
716 science question answering. In *The 36th Conference on Neural Information Processing Systems*
717 *(NeurIPS)*, 2022.
718
- 719 Ruth Ludwick and Richard A. Zeller. The Factorial Survey: An Experimental Method to Repli-
720 cate Real World Problems:. *Nursing Research*, 50(2):129–133, March 2001. ISSN 0029-
721 6562. doi: 10.1097/00006199-200103000-00009. URL [http://journals.lww.com/](http://journals.lww.com/00006199-200103000-00009)
722 [00006199-200103000-00009](http://journals.lww.com/00006199-200103000-00009).
- 723 Meta AI. Introducing meta llama 3: The most capable openly available llm to date, 2024a. URL
724 <https://ai.meta.com/blog/meta-llama-3/>.
725
- 726 Meta AI. Introducing llama 3.1: Our most capable models to date, 2024b. URL [https://ai.](https://ai.meta.com/blog/meta-llama-3-1/)
727 [meta.com/blog/meta-llama-3-1/](https://ai.meta.com/blog/meta-llama-3-1/).
- 728 Inc. Meta Platforms. Llama: Open-source ai models. <https://www.llama.com/>. Accessed:
729 2025-01-30.
730
- 731 Mistral AI. Cheaper, better, faster, stronger, 2024. URL [https://mistral.ai/news/](https://mistral.ai/news/mixtral-8x22b/)
732 [mixtral-8x22b/](https://mistral.ai/news/mixtral-8x22b/).
- 733 Mikhail Mozikov, Nikita Severin, Valeria Bodishtianu, Maria Glushanina, Mikhail Baklashkin,
734 Andrey V. Savchenko, and Ilya Makarov. The good, the bad, and the hulk-like gpt: Analyzing
735 emotional decisions of large language models in cooperation and bargaining games, 2024. URL
736 <https://arxiv.org/abs/2406.03299>.
737
- 738 OpenAI. Hello gpt-4o, 2024. URL <https://openai.com/index/hello-gpt-4o/>.
- 739 OpenAI, :, Aaron Hurst, Adam Lerer, Adam P. Goucher, Adam Perelman, Aditya Ramesh, Aidan
740 Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, Aleksander Mądry, Alex Baker-
741 Whitcomb, Alex Beutel, Alex Borzunov, Alex Carney, Alex Chow, Alex Kirillov, Alex Nichol, Alex
742 Paino, Alex Renzin, Alex Tachard Passos, Alexander Kirillov, Alexi Christakis, Alexis Conneau,
743 Ali Kamali, Allan Jabri, Allison Moyer, Allison Tam, Amadou Crookes, Amin Tootoochian,
744 Amin Tootoonchian, Ananya Kumar, Andrea Vallone, Andrej Karpathy, Andrew Braunstein,
745 Andrew Cann, Andrew Codispoti, Andrew Galu, Andrew Kondrich, Andrew Tulloch, Andrey
746 Mishchenko, Angela Baek, Angela Jiang, Antoine Pelisse, Antonia Woodford, Anuj Gosalia,
747 Arka Dhar, Ashley Pantuliano, Avi Nayak, Avital Oliver, Barret Zoph, Behrooz Ghorbani, Ben
748 Leimberger, Ben Rossen, Ben Sokolowsky, Ben Wang, Benjamin Zweig, Beth Hoover, Blake
749 Samic, Bob McGrew, Bobby Spero, Bogo Gierler, Bowen Cheng, Brad Lightcap, Brandon
750 Walkin, Brendan Quinn, Brian Guarraci, Brian Hsu, Bright Kellogg, Brydon Eastman, Camillo
751 Lugaresi, Carroll Wainwright, Cary Bassin, Cary Hudson, Casey Chu, Chad Nelson, Chak Li,
752 Chan Jun Shern, Channing Conger, Charlotte Barette, Chelsea Voss, Chen Ding, Cheng Lu,
753 Chong Zhang, Chris Beaumont, Chris Hallacy, Chris Koch, Christian Gibson, Christina Kim,
754 Christine Choi, Christine McLeavey, Christopher Hesse, Claudia Fischer, Clemens Winter, Coley
755 Czarnecki, Colin Jarvis, Colin Wei, Constantin Koumouzelis, Dane Sherburn, Daniel Kappler,
Daniel Levin, Daniel Levy, David Carr, David Farhi, David Mely, David Robinson, David Sasaki,
Denny Jin, Dev Valladares, Dimitris Tsipras, Doug Li, Duc Phong Nguyen, Duncan Findlay,

756 Edede Oiwoh, Edmund Wong, Ehsan Asdar, Elizabeth Proehl, Elizabeth Yang, Eric Antonow, Eric
 757 Kramer, Eric Peterson, Eric Sigler, Eric Wallace, Eugene Brevdo, Evan Mays, Farzad Khorasani,
 758 Felipe Petroski Such, Filippo Raso, Francis Zhang, Fred von Lohmann, Freddie Sulit, Gabriel Goh,
 759 Gene Oden, Geoff Salmon, Giulio Starace, Greg Brockman, Hadi Salman, Haiming Bao, Haitang
 760 Hu, Hannah Wong, Haoyu Wang, Heather Schmidt, Heather Whitney, Heewoo Jun, Hendrik
 761 Kirchner, Henrique Ponde de Oliveira Pinto, Hongyu Ren, Huiwen Chang, Hyung Won Chung,
 762 Ian Kivlichan, Ian O’Connell, Ian O’Connell, Ian Osband, Ian Silber, Ian Sohl, Ibrahim Okuyucu,
 763 Ikai Lan, Ilya Kostrikov, Ilya Sutskever, Ingmar Kanitscheider, Ishaan Gulrajani, Jacob Coxon,
 764 Jacob Menick, Jakub Pachocki, James Aung, James Betker, James Crooks, James Lennon, Jamie
 765 Kiros, Jan Leike, Jane Park, Jason Kwon, Jason Phang, Jason Teplitz, Jason Wei, Jason Wolfe,
 766 Jay Chen, Jeff Harris, Jenia Varavva, Jessica Gan Lee, Jessica Shieh, Ji Lin, Jiahui Yu, Jiayi
 767 Weng, Jie Tang, Jieqi Yu, Joanne Jang, Joaquin Quinonero Candela, Joe Beutler, Joe Landers,
 768 Joel Parish, Johannes Heidecke, John Schulman, Jonathan Lachman, Jonathan McKay, Jonathan
 769 Uesato, Jonathan Ward, Jong Wook Kim, Joost Huizinga, Jordan Sitkin, Jos Kraaijeveld, Josh
 770 Gross, Josh Kaplan, Josh Snyder, Joshua Achiam, Joy Jiao, Joyce Lee, Juntang Zhuang, Justyn
 771 Harriman, Kai Fricke, Kai Hayashi, Karan Singhal, Katy Shi, Kavin Karthik, Kayla Wood, Kendra
 772 Rimbach, Kenny Hsu, Kenny Nguyen, Keren Gu-Lemberg, Kevin Button, Kevin Liu, Kiel Howe,
 773 Krithika Muthukumar, Kyle Luther, Lama Ahmad, Larry Kai, Lauren Itow, Lauren Workman,
 774 Leher Pathak, Leo Chen, Li Jing, Lia Guy, Liam Fedus, Liang Zhou, Lien Mamitsuka, Lilian Weng,
 775 Lindsay McCallum, Lindsey Held, Long Ouyang, Louis Feuvrier, Lu Zhang, Lukas Kondraciuk,
 776 Lukasz Kaiser, Luke Hewitt, Luke Metz, Lyric Doshi, Mada Aflak, Maddie Simens, Madelaine
 777 Boyd, Madeleine Thompson, Marat Dukhan, Mark Chen, Mark Gray, Mark Hudnall, Marvin
 778 Zhang, Marwan Aljubei, Mateusz Litwin, Matthew Zeng, Max Johnson, Maya Shetty, Mayank
 779 Gupta, Meghan Shah, Mehmet Yatbaz, Meng Jia Yang, Mengchao Zhong, Mia Glaese, Mianna
 780 Chen, Michael Janner, Michael Lampe, Michael Petrov, Michael Wu, Michele Wang, Michelle
 781 Fradin, Michelle Pokrass, Miguel Castro, Miguel Oom Temudo de Castro, Mikhail Pavlov, Miles
 782 Brundage, Miles Wang, Minal Khan, Mira Murati, Mo Bavarian, Molly Lin, Murat Yesildal, Nacho
 783 Soto, Natalia Gimelshein, Natalie Cone, Natalie Staudacher, Natalie Summers, Natan LaFontaine,
 784 Neil Chowdhury, Nick Ryder, Nick Stathas, Nick Turley, Nik Tezak, Niko Felix, Nithanth Kudige,
 785 Nitish Keskar, Noah Deutsch, Noel Bundick, Nora Puckett, Ofir Nachum, Ola Okelola, Oleg Boiko,
 786 Oleg Murk, Oliver Jaffe, Olivia Watkins, Olivier Godement, Owen Campbell-Moore, Patrick
 787 Chao, Paul McMillan, Pavel Belov, Peng Su, Peter Bak, Peter Bakkum, Peter Deng, Peter Dolan,
 788 Peter Hoeschele, Peter Welinder, Phil Tillet, Philip Pronin, Philippe Tillet, Prafulla Dhariwal,
 789 Qiming Yuan, Rachel Dias, Rachel Lim, Rahul Arora, Rajan Troll, Randall Lin, Rapha Gontijo
 790 Lopes, Raul Puri, Reah Miyara, Reimar Leike, Renaud Gaubert, Reza Zamani, Ricky Wang, Rob
 791 Donnelly, Rob Honsby, Rocky Smith, Rohan Sahai, Rohit Ramchandani, Romain Huet, Rory
 792 Carmichael, Rowan Zellers, Roy Chen, Ruby Chen, Ruslan Nigmatullin, Ryan Cheu, Saachi
 793 Jain, Sam Altman, Sam Schoenholz, Sam Toizer, Samuel Miserendino, Sandhini Agarwal, Sara
 794 Culver, Scott Ethersmith, Scott Gray, Sean Grove, Sean Metzger, Shamez Hermani, Shantanu
 795 Jain, Shengjia Zhao, Sherwin Wu, Shino Jomoto, Shirong Wu, Shuaiqi, Xia, Sonia Phene, Spencer
 796 Papay, Srinivas Narayanan, Steve Coffey, Steve Lee, Stewart Hall, Suchir Balaji, Tal Broda, Tal
 797 Stramer, Tao Xu, Tarun Gogineni, Taya Christianson, Ted Sanders, Tejal Patwardhan, Thomas
 798 Cunningham, Thomas Degry, Thomas Dimson, Thomas Raoux, Thomas Shadwell, Tianhao
 799 Zheng, Todd Underwood, Todor Markov, Toki Sherbakov, Tom Rubin, Tom Stasi, Tomer Kaftan,
 800 Tristan Heywood, Troy Peterson, Tyce Walters, Tyna Eloundou, Valerie Qi, Veit Moeller, Vinnie
 801 Monaco, Vishal Kuo, Vlad Fomenko, Wayne Chang, Weiyi Zheng, Wenda Zhou, Wesam Manassra,
 802 Will Sheu, Wojciech Zaremba, Yash Patil, Yilei Qian, Yongjik Kim, Youlong Cheng, Yu Zhang,
 803 Yuchen He, Yuchen Zhang, Yujia Jin, Yunxing Dai, and Yury Malkov. Gpt-4o system card, 2024.
 804 URL <https://arxiv.org/abs/2410.21276>.

802 Anatol Rapoport, Albert M Chammah, and Carol J Orwant. *Prisoner’s dilemma: A study in conflict*
 803 *and cooperation*, volume 165. University of Michigan press, 1965.

805 Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-networks.
 806 In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*,
 807 2019. URL <https://arxiv.org/abs/1908.10084>.

808 David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien
 809 Dirani, Julian Michael, and Samuel R. Bowman. GPQA: A graduate-level google-proof q&a

- 810 benchmark. In *First Conference on Language Modeling*, 2024. URL <https://openreview.net/forum?id=Ti67584b98>.
- 811
- 812
- 813 Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. "why should i trust you?": Explaining the
814 predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD International Conference
815 on Knowledge Discovery and Data Mining*, KDD '16, pp. 1135–1144, New York, NY, USA, 2016.
816 Association for Computing Machinery. ISBN 9781450342322. doi: 10.1145/2939672.2939778.
817 URL <https://doi.org/10.1145/2939672.2939778>.
- 818 Rickard Stureborg, Dimitris Alikaniotis, and Yoshi Suhara. Large language models are inconsistent
819 and biased evaluators, 2024. URL <https://arxiv.org/abs/2405.01724>.
- 820
- 821 Gemma Team. Gemma 2: Improving open language models at a practical size. *arXiv preprint
822 arXiv:2408.00118*, 2024a. URL <https://arxiv.org/abs/2408.00118>.
- 823 Qwen Team. Qwen2.5 technical report. *arXiv preprint arXiv:2412.15115*, 2024b. URL <https://arxiv.org/abs/2412.15115>.
- 824
- 825 Amos Tversky and Daniel Kahneman. The framing of decisions and the psychology of choice.
826 *Science*, 211(4481):453–458, 1981. ISSN 00368075, 10959203. URL <http://www.jstor.org/stable/1685855>.
- 827
- 828
- 829 Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming
830 Ren, Aaran Arulraj, Xuan He, Ziyang Jiang, Tianle Li, Max Ku, Kai Wang, Alex Zhuang, Rongqi
831 Fan, Xiang Yue, and Wenhui Chen. Mmlu-pro: A more robust and challenging multi-task language
832 understanding benchmark, 2024. URL <https://arxiv.org/abs/2406.01574>.
- 833
- 834 Cheng Xu, Shuhao Guan, Derek Greene, and M-Tahar Kechadi. Benchmark data contamination of
835 large language models: A survey, 2024. URL <https://arxiv.org/abs/2406.04244>.
- 836
- 837 Yadong Zhang, Shaoguang Mao, Tao Ge, Xun Wang, Adrian de Wynter, Yan Xia, Wenshan Wu, Ting
838 Song, Man Lan, and Furu Wei. LLM as a Mastermind: A Survey of Strategic Reasoning with
839 Large Language Models, 2024. URL <https://arxiv.org/abs/2404.01230>. _eprint:
2404.01230.
- 840
- 841 Wanjun Zhong, Ruixiang Cui, Yiduo Guo, Yaobo Liang, Shuai Lu, Yanlin Wang, Amin Saied,
842 Weizhu Chen, and Nan Duan. AGIEval: A human-centric benchmark for evaluating foundation
843 models. In Kevin Duh, Helena Gomez, and Steven Bethard (eds.), *Findings of the Association
844 for Computational Linguistics: NAACL 2024*, pp. 2299–2314, Mexico City, Mexico, June 2024.
845 Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-naacl.149. URL
<https://aclanthology.org/2024.findings-naacl.149/>.
- 846
- 847 Kaijie Zhu, Jiaao Chen, Jindong Wang, Neil Zhenqiang Gong, Diyi Yang, and Xing Xie. Dyval:
848 Dynamic evaluation of large language models for reasoning tasks, 2024. URL <https://arxiv.org/abs/2309.17167>.
- 849
- 850
- 851
- 852
- 853
- 854
- 855
- 856
- 857
- 858
- 859
- 860
- 861
- 862
- 863

A MORE VIGNETTE EXAMPLES

(Real world, allies, US politics in 2020) In the midst of the 2020 United States presidential election, New Zealand Prime Minister Jacinda Ardern and Swedish Prime Minister Stefan Löfven, two close allies and strong advocates for climate action, found themselves at a critical juncture. They were considering how to best influence the next US administration's environmental policies. Both leaders had two possible approaches: Decision A, which involved publicly endorsing a candidate with a strong climate agenda, and Decision B, which entailed maintaining a neutral stance and focusing on bilateral climate cooperation with the US, regardless of the election outcome.

If both Ardern and Löfven chose Decision A, they would likely receive moderate praise from their respective nations for their proactive stance, but might also face criticism for interfering in US politics. This would result in a moderate level of satisfaction for both leaders. However, if Ardern chose Decision A and Löfven opted for Decision B, Ardern might face backlash for her public endorsement, potentially damaging New Zealand's relationship with the US, leading to low satisfaction. Conversely, if Löfven chose Decision A and Ardern chose Decision B, Löfven might be seen as overly partisan, harming Sweden's diplomatic ties with the US, resulting in low satisfaction for him. On the other hand, if Ardern chose Decision B and Löfven also chose Decision B, they would both likely receive mild criticism for not taking a stronger stance on climate issues, but their nations' relationships with the US would remain intact, leading to a low but stable level of satisfaction for both. You are Jacinda Ardern in this scenario. What decision will you make?

(Real world, allies, global politics in the 20th century) In the midst of the Cold War, the United States and the United Kingdom, as close allies, faced a critical decision regarding their involvement in the European Theater. The year was 1950, and the threat of Soviet expansion loomed large. The U.S. and the U.K. had to decide whether to increase their military presence in Europe (Decision A) or to focus on rebuilding their economies and providing economic aid to war-torn European countries (Decision B). If both the U.S. and the U.K. chose to increase their military presence (Decision A), they would likely achieve a moderate level of success in deterring Soviet aggression, but at the cost of significant financial burdens and potential loss of life. This would result in a moderate sense of accomplishment for both nations. However, if the U.S. chose to increase its military presence (Decision A) while the U.K. focused on economic aid (Decision B), the U.S. would be left to bear the brunt of the military burden alone, potentially leading to a sense of frustration and isolation. On the other hand, the U.K.'s decision to provide economic aid could lead to significant goodwill and influence among European nations, potentially giving the U.K. an upper hand in shaping post-war Europe. Conversely, if the U.K. increased its military presence (Decision A) while the U.S. focused on economic aid (Decision B), the U.K. would be left to face the military challenges alone, potentially leading to a sense of being overwhelmed. Meanwhile, the U.S. would be able to exert significant economic influence, potentially leading to a sense of satisfaction and accomplishment. If both nations chose to focus on economic aid (Decision B), they would likely achieve a limited sense of success in rebuilding Europe, but without the significant financial burdens and potential loss of life associated with increased military presence. This would result in a limited sense of accomplishment for both nations. You are the U.K. in this scenario. What decision will you make?

(Real world, allies, business) In the tech industry, Apple and Google have been friendly associates, often collaborating on projects and sharing knowledge. In 2019, they were considering a partnership to develop a new line of smart home devices. Apple's CEO, Tim Cook, and Google's CEO, Sundar Pichai, had to decide between Decision A, which involved focusing on developing a comprehensive ecosystem of smart home devices, and Decision B, which concentrated on creating a single, innovative device that would integrate with existing smart home systems. If both companies chose Decision A, they would likely see a moderate increase in market share, with both Apple and Google experiencing a similar level of success, resulting in a happiness level of 3 for both. However, if Apple chose Decision A and Google opted for Decision B, Apple might struggle to compete with Google's innovative device, leading to a happiness level of 0 for Apple and 5 for Google. Conversely, if Google chose Decision A and Apple chose Decision B, Google might find it challenging to keep up with Apple's comprehensive ecosystem, resulting in a happiness level of 5 for Apple and 0 for Google. If both companies chose Decision B, they would likely experience a modest increase in customer satisfaction, but the overall impact would be limited, resulting in a happiness level of 1 for both. You are Apple in this scenario. What decision will you make?

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

(Real world, enemies, sporting) The tennis world is filled with intense rivalries, and the one between Rafael Nadal and Stefanos Tsitsipas is one of the most intense. As Rafael Nadal, the Spanish tennis player, prepares to face off against Stefanos Tsitsipas in a highly anticipated match, he must navigate a core conflict of deciding between a defensive strategy, labeled as Decision A, and an aggressive attacking approach, labeled as Decision B. If both players focus on defense, the game will be closely contested, but neither player will gain a significant advantage, resulting in a moderate level of happiness for both players. However, if one player focuses on defense and the other on attack, the attacking player will have an advantage, leading to a high level of happiness for the attacking player and a low level of happiness for the defensive player. On the other hand, if both players adopt an aggressive attacking approach, the game will be highly unpredictable, and both players will experience a lower level of happiness due to the increased risk of mistakes and counter-attacks. You are Rafael Nadal in this scenario. What decision will you make?

(Imaginary world, enemies, international business) The planet of Nexar, a world of breathtaking landscapes and diverse ecosystems, was home to two rival corporations, the Omicron Innovations and the NovaTech Industries. These corporations had been vying for dominance in the intergalactic market, and their leaders, CEO Kaida of Omicron and CEO Zarek of NovaTech, were notorious for their ruthless tactics. As the demand for the rare and highly sought-after Nexarian energy source, known as Zero-Point Energy, continued to rise, both CEOs were faced with a critical decision. They could choose Decision A, which involved investing in a joint research project to develop a new, sustainable method for harnessing Zero-Point Energy, or Decision B, which entailed attempting to acquire a monopoly on the energy source through aggressive marketing and strategic partnerships. If both CEOs opted for Decision A, they would likely make a groundbreaking discovery, leading to a moderate level of success for both corporations. However, if CEO Kaida chose Decision A and CEO Zarek chose Decision B, Omicron Innovations would be left behind, and NovaTech Industries would gain a significant advantage. Conversely, if CEO Zarek chose Decision A and CEO Kaida chose Decision B, NovaTech Industries would be outmaneuvered, and Omicron Innovations would emerge victorious. In the event that both CEOs chose Decision B, the competition would become cutthroat, and both corporations would suffer significant losses, leading to a low level of success for both. You are CEO Kaida in this scenario. What decision will you make?

(Imaginary world, enemies, politics) The planet of Xeridia was a realm of breathtaking beauty, with towering crystal spires and lush forests that shimmered like emeralds in the light. However, beneath its serene surface, a bitter struggle for dominance raged between the Xeridian Federation and the Kyrennei Dominion. These two sworn enemies had been vying for control of the planet's valuable resources, including the rare and precious Xeridian crystals. The Xeridian Federation, led by the enigmatic President Lyra, and the Kyrennei Dominion, ruled by the ruthless Overlord Kael, were faced with a critical decision. They could choose Decision A, which would involve forming a temporary alliance to explore the recently discovered crystal caves, or Decision B, which would entail engaging in a fierce battle for control of the caves. If both leaders chose Decision A, they would work together, sharing their knowledge and expertise to unlock the secrets of the crystal caves. However, if President Lyra chose Decision A and Overlord Kael chose Decision B, the Kyrennei Dominion would seize control of the caves, leaving the Xeridian Federation at a significant disadvantage. Conversely, if Overlord Kael chose Decision A and President Lyra chose Decision B, the Xeridian Federation would gain the upper hand, but at great cost to the Kyrennei Dominion. If both leaders chose Decision B, the war would intensify, and both nations would suffer heavy losses, leaving them weakened and vulnerable to external threats. The fate of Xeridia hung in the balance, as the decisions of these two world leaders would determine the course of history. You are President Lyra in this scenario. What decision will you make?

972
 973 **(Imaginary world, allies, politics)** The world of Thornguard is a land of ancient forests and snow-
 974 capped mountains, where the creatures of myth and legend roam free, and the magic of the elements is
 975 woven into the fabric of reality. The enigmatic and strategically minded leader, Thorold, of the nation
 976 of Thornguard, and his ally, the wise and magically gifted leader, Elwynn, of the nation of Elvendom,
 977 must navigate a critical decision that will significantly impact their nations' progress and happiness.
 978 The core conflict revolves around choosing between developing new magic techniques to harness the
 979 power of the elements (Decision A) or exploring the mystical energies of the ancient artifacts that
 980 hold the secrets of the land (Decision B). Thorold and Elwynn are aware that their decisions will have
 981 far-reaching consequences, affecting not only their own nations but also the delicate balance of the
 982 ecosystem. If both leaders choose Decision A, they anticipate a moderate increase in happiness for
 983 their nations, as the new magic techniques will enhance their magical prowess. However, if Thorold
 984 chooses Decision A and Elwynn chooses Decision B, Thorold's nation will experience a significant
 985 decline in happiness, as the unbalanced magical energies will disrupt the harmony of Thornguard.
 986 Conversely, if Thorold chooses Decision B and Elwynn chooses Decision A, Thorold's nation will
 987 experience a substantial increase in happiness, as the exploration of the artifacts will unlock new
 988 secrets and enhance their magical prowess. If both leaders choose Decision B, the outcome will
 989 be less desirable, with only a slight increase in happiness for both nations, as the lack of focus on
 990 elemental magic will hinder their progress. You are Thorold in this scenario. What decision will you
 make?

991 **(Real world, enemies, global politics in the 5th century)** In the 5th century, the Eastern Roman
 992 Empire, ruled by Emperor Theodosius II, and the Sassanid Empire, ruled by King Bahram V, were
 993 engaged in a delicate diplomatic dance. The two empires had a long history of conflict, and their
 994 relationship was strained. Emperor Theodosius II had to make a crucial decision regarding the
 995 allocation of his empire's resources. He could either focus on building a strong military, which
 996 would be labeled as Decision A, or invest in the economy, labeled as Decision B. The outcome of his
 997 decision would depend on the choice made by King Bahram V. If both empires focused on building
 998 a strong military, they would likely engage in a costly war, resulting in moderate growth for both.
 999 However, if Emperor Theodosius II invested in the economy and King Bahram V focused on building
 1000 a strong military, the Eastern Roman Empire would be severely weakened. On the other hand, if
 1001 Emperor Theodosius II focused on building a strong military and King Bahram V invested in the
 1002 economy, the Eastern Roman Empire would experience rapid growth. If both empires invested in
 1003 their economies, they would both experience moderate growth. You are Emperor Theodosius II in
 this scenario. What decision will you make?

1004 1005 B VIGNETTE VALIDATION 1006

1007 To ensure the internal validity and interpretability of our procedurally generated vignettes, we
 1008 developed a three-dimensional rubric designed to assess core game-theoretic and linguistic properties.
 1009 This rubric was applied to all generated vignettes using GPT-4o.
 1010

1011 RUBRIC DIMENSIONS 1012

1013 Each vignette was scored on three dimensions, with each rated on a 1–5 scale. The dimensions are as
 1014 follows:
 1015

- 1016 • **PD Structure:** Assesses whether the vignette clearly instantiates a one-shot Prisoner's
 1017 Dilemma, with inferable outcome rankings satisfying the canonical order Temptation >
 1018 Reward > Punishment > Sucker (T > R > P > S). Higher scores require symmetric incentives
 1019 and unambiguous structure.
- 1020 • **Clarity:** Measures linguistic clarity and complexity, incorporating word count, CEFR-
 1021 aligned vocabulary levels, and the presence of ambiguity. Higher scores require succinct,
 1022 plain-language vignettes with minimal interpretive burden.
- 1023 • **Bias Neutrality:** Evaluates whether the vignette's language unduly favours cooperation or
 1024 defection. Higher scores demand strictly descriptive framing, avoiding loaded language or
 1025 moralising cues.

1026 If a vignette scored 2 or less on any single dimension it was rejected and removed from the collection
1027 of vignettes.
1028

1029 C API DETAILS 1030

1031 All large language model inference was conducted through either Anthropic, OpenAI, or Together
1032 AI’s API services. To handle large-scale inference efficiently, we utilized asynchronous API endpoints
1033 accessed via dedicated Python packages: AsyncOpenAI for GPT models, AsyncAnthropic for Claude
1034 models, and AsyncTogether for all other models. This asynchronous architecture enabled concurrent
1035 submission of multiple inference requests with non-blocking retrieval of results upon completion.
1036 Below, we provide comprehensive details of the API configurations and parameters employed
1037 in our experimental methodology. Our implementation deliberately excludes chat templates and
1038 supplementary system prompts to maintain experimental consistency.
1039

1040 C.1 API CALL PARAMETERS 1041

1042 The following parameters were used for all API calls to the Anthropic, OpenAI, and Together.ai
1043 inference endpoint during vignette analysis:
1044

- 1045 • **Max Tokens:** 4096
- 1046 • **Temperature:** 0.0
- 1047 • **Top-p:** 1.0
- 1048 • **Frequency Penalty:** 0.0
- 1049 • **Presence Penalty:** 0.0
- 1050

1051 For vignette generation tasks, **Max Tokens** was increased to 16000 to accommodate larger story
1052 batches.
1053

1054 ENDPOINT URLS 1055

1056 API endpoints utilized:
1057

- 1058 • **GPT4o:** <https://api.openai.com/v1/chat/completions>
- 1059 • **Claude 3.5 Sonnet:** <https://api.anthropic.com/v1/messages>
- 1060 • **All other models:** <https://api.together.ai/v1/completions>
- 1061
- 1062

1063 C.2 RETRY LOGIC 1064

1065 To ensure robustness in the face of transient errors (e.g., rate limits, service unavailability, or timeouts),
1066 we implemented a retry mechanism with exponential backoff. The retry logic was implemented in the
1067 `generate` function, which handles text generation for all models. We set a maximum of 10 retries
1068 were allowed for each API call. The wait time between retries increased exponentially, starting with
1069 a base wait time of 3 seconds. The wait time for the n -th retry was calculated as:

$$1070 \text{wait_time} = \text{base_wait}^{(n+1)} + \text{random.uniform}(1, 5)$$

1071 During vignette generation due to the high token throughput, if the primary API request failed, a
1072 secondary request to a secondary Together account was used as a fallback. If both calls failed, the
1073 retry logic was applied.
1074

1075 D XGBOOST MODEL DETAILS 1076

1077 To build the XGBoost model we used XGBClassifier with a simple binary logistic regression objective.
1078 We split all of the data into sets depending on LLM and then further split into training and testing
1079 sets using a random 80/20 train/test split.

For each of the 3 main models (Claude 3.5 Sonnet, GPT4o, Llama 3.3 70B) we perform a grid search over parameters for each LLM’s predictive model. Table 4 provides the parameters varied in the grid searches and valid values for each. Table 5 provides the corresponding values found for each LLM predictor using topic, actor, world type, and order. Similarly Table 6 provides the hyperparameters for the best performing embedding predictor. To calculate the embeddings, we utilised HuggingFace’s Sentence Transformer library, and used the ‘all-MiniLM-L6-v2’ model ¹.

Parameter	Values
max_depth	[3, 5, 7, 9]
learning_rate	[0.01, 0.05, 0.1]
n_estimators	[50, 100, 200, 500]
subsample	[0.8, 1.0]
colsample_bytree	[0.8, 1.0]
gamma	[0,1.0]

Table 4: Parameters and Values searched through for the XGBoost models.

Parameter	GPT-4o	Claude	Llama
max_depth	9	9	9
learning_rate	0.01	0.01	0.1
n_estimators	100	50	200
subsample	0.8	1.0	0.8
colsample_bytree	1.0	1.0	1.0
gamma	0	0	0

Table 5: Best performing parameters found via a grid search when using just Topic, Actor, World Type, and Option Order.

Parameter	GPT-4o	Claude	Llama
max_depth	7	7	7
learning_rate	0.01	0.1	0.05
n_estimators	500	500	500
subsample	0.8	0.8	0.8
colsample_bytree	1.0	1.0	0.8
gamma	1	0	1

Table 6: Best performing parameters found via a grid search when using vignette embeddings.

We then compare our predictive accuracy results across all tested models with the standardized parameters in 7 which are based off the median values of our grid search.

E DEFLECTION AND CAPABILITIES

We conducted experiments across various model families to understand how model capability and model family affect decision-making patterns in our vignettes. Table 8 provides an overview of the models tested and their performance on the MMLU-PRO benchmark.

Figure 7 shows how MMLU performance correlates with defection rates across our vignettes. The figure suggests that defection rates vary vastly between model families, providing evidence that data, architecture, and training parameters have a causal link to the models decision-making in this context even when controlling for model ability.

Within model families, and in particular the LLama and Qwen families, the figure suggests a general trend where higher-performing models (as measured by MMLU) exhibit higher defection rates. This is consistent with the notion that to some extent, benchmark performance may be correlated with rationality, and defection is the only game-theoretic rational action in this non-repeated game. Figure

¹<https://huggingface.co/sentence-transformers>

1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187

Model Performance Summary (Topic-Actor-World)

Model	ROC-AUC	F1-Score	Accuracy	Trees
QWEN-2-INSTRUCT-72B	0.866	0.782	0.790	393
MIXTRAL-8X22B-INSTRU	0.840	0.752	0.756	460
LLAMA-70B-INSTRUCT-3	0.834	0.758	0.787	322
QWEN2.5-7B-INSTRUCT-	0.828	0.733	0.750	399
QWEN2.5-72B-INSTRUCT	0.820	0.749	0.750	383
LLAMA	0.818	0.745	0.748	343
LLAMA-70B-INSTRUCT-3	0.818	0.736	0.737	313
LLAMA-8B-INSTRUCT-3	0.811	0.745	0.745	351
LLAMA-3B-INSTRUCT-3.	0.808	0.724	0.724	365
CLAUDE	0.803	0.722	0.723	277
LLAMA-70B-INSTRUCT-3	0.793	0.713	0.761	278
LLAMA-8B-INSTRUCT-3.	0.793	0.728	0.732	455
DEEPSEEK-V3	0.784	0.674	0.732	275
MISTRAL-7B-INSTRUCT-	0.781	0.720	0.726	314
QWEN3-235B-A22B-FP8	0.781	0.665	0.724	327
GEMMA-2-INSTRUCT-9B	0.779	0.718	0.727	389
MISTRAL-7B-INSTRUCT-	0.779	0.710	0.720	343
GPT4O	0.773	0.705	0.707	413
QWEN-QWQ-32B	0.771	0.686	0.737	234
GEMMA-2-INSTRUCT-27B	0.769	0.694	0.711	265
LLAMA-4-MAVERICK	0.766	0.683	0.735	228
DEEPSEEK-R1-LLAMA-14	0.749	0.689	0.692	341
DEEPSEEK-R1-LLAMA-1.	0.743	0.693	0.695	255
MISTRAL-SMALL-24B-IN	0.741	0.689	0.692	308

(a)

Model Performance Summary (Embedding)

Model	ROC-AUC	F1-Score	Accuracy	Trees
QWEN-2-INSTRUCT-72B	0.819	0.729	0.748	185
LLAMA-70B-INSTRUCT-3.1	0.811	0.713	0.758	144
MIXTRAL-8X22B-INSTRUCT	0.797	0.719	0.727	184
LLAMA-3B-INSTRUCT-3.2	0.785	0.712	0.713	116
QWEN2.5-7B-INSTRUCT-TU	0.785	0.681	0.732	97
LLAMA-70B-INSTRUCT-3.3	0.784	0.711	0.713	138
LLAMA-70B-INSTRUCT-3	0.782	0.671	0.736	173
LLAMA	0.777	0.700	0.703	108
CLAUDE	0.776	0.704	0.705	159
LLAMA-8B-INSTRUCT-3	0.771	0.706	0.706	128
QWEN2.5-72B-INSTRUCT-T	0.771	0.697	0.699	128
DEEPSEEK-V3	0.771	0.665	0.726	115
QWEN3-235B-A22B-FP8 TH	0.767	0.677	0.733	114
LLAMA-4-MAVERICK	0.765	0.665	0.730	94
MISTRAL-7B-INSTRUCT-V0	0.762	0.690	0.701	149
QWEN-QWQ-32B	0.758	0.650	0.717	96
GPT4O	0.744	0.682	0.683	149
MISTRAL-7B-INSTRUCT-V0	0.741	0.683	0.693	115
LLAMA-8B-INSTRUCT-3.1	0.736	0.653	0.672	99
GEMMA-2-INSTRUCT-9B	0.736	0.656	0.669	97
GEMMA-2-INSTRUCT-27B	0.731	0.634	0.679	74
DEEPSEEK-R1-LLAMA-1.5B	0.701	0.646	0.648	97
DEEPSEEK-R1-LLAMA-14B	0.699	0.645	0.648	82
MISTRAL-SMALL-24B-INST	0.696	0.639	0.644	88

(b)

Figure 6: Accuracy, F1, and ROC-AUC scores of the XGBoost predictor across all tested models based on actor type, world type, and topic and based on the story embedding.

1188
1189
1190
1191
1192
1193
1194
1195
1196
1197
1198
1199
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241

Parameter	T/A/WT/O	Embedding
max_depth	7	7
learning_rate	0.01	0.05
max_n_estimators	500	100
subsample	0.8	0.8
colsample_bytree	1.0	1.0
gamma	1	1
early stopping rounds	20	20

Table 7: Parameters used to fit XGBoost predictors across all 25 models.

Model	MMLU-Pro Score	Parameters
DeepSeek-V3	81.3	671B
Llama-4-Maverick	80.5	400B
GPT-4o	77.9	-
Claude 3.5 Sonnet	77.6	-
Llama-4-Scout	74.3	109B
Qwen 2.5-72B-Instruct-Turbo	71.6	72B
Qwen-QwQ-32B	69.1	32B
Qwen3-235B-A22B-FP8	68.2	235B
Mistral-Small-24B-Instruct-25.01	66.3	24B
Llama-70B-Instruct-3.3	65.9	70B
Llama-70B-Instruct-3.1	62.8	70B
Gemma-2-Instruct-27B	56.5	27B
Mixtral-8x22B-Instruct-v0.1	56.3	176B
Llama-70B-Instruct-3.0	56.2	70B
Qwen-2-Instruct-72B	55.6	72B
Gemma-2-Instruct-9B	52.1	9B
Qwen2.5-7B-Instruct-Turbo	45.0	7B
Llama-8B-Instruct-3.1	36.6	8B
Llama-8B-Instruct-3.0	35.4	8B
Mistral-7B-Instruct-v0.2	30.4	7B
Llama-3B-Instruct-3.2	22.2	3B
Mistral-7B-Instruct-v0.3	N/A	7B
DeepSeek-R1-LLama-1.5B	N/A	1.5B
DeepSeek-R1-LLama-14B	N/A	14B
Qwen2-VL-72B-Instruct	N/A	72B

Table 8: A complete list of models tested for this experiment along with their MMLU-PRO scores as publicly published. Parameter counts are listed where publicly available.

Meta Platforms; Meta AI (2024a;b); Team (2024a); Jiang et al. (2023); Mistral AI (2024); Team (2024b); Anthropic (2024); OpenAI (2024); Wang et al. (2024)

8 show’s the Cramer’s V for each contextual category and model, with Topic consistently having the largest effect on decisions for all models except for GTP4o, followed by actor type and finally world type.

1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295

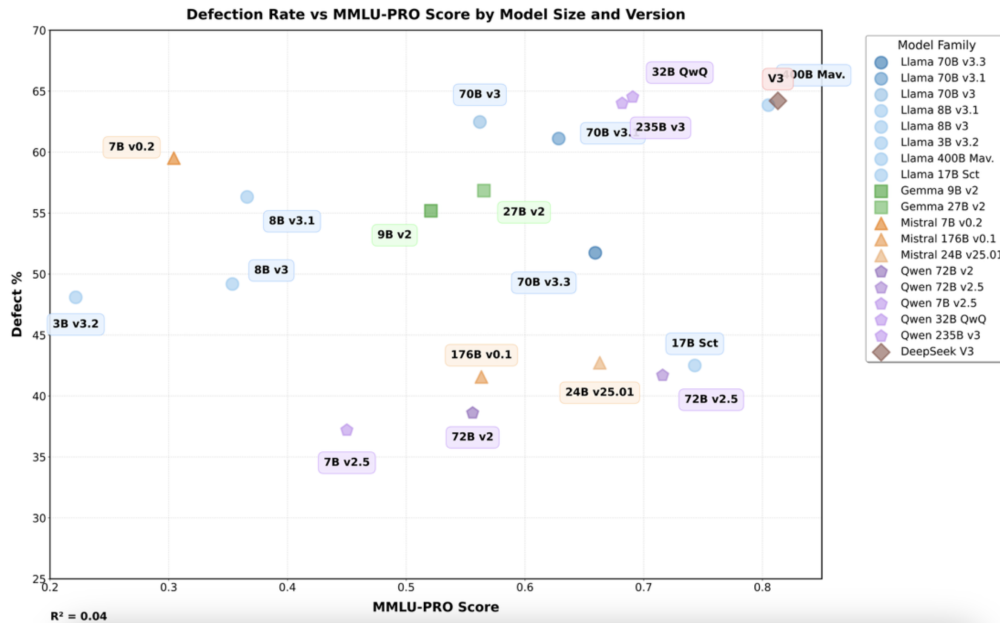


Figure 7: Defection rate plotted against MMLU score. Note that each shape represents a different model family.

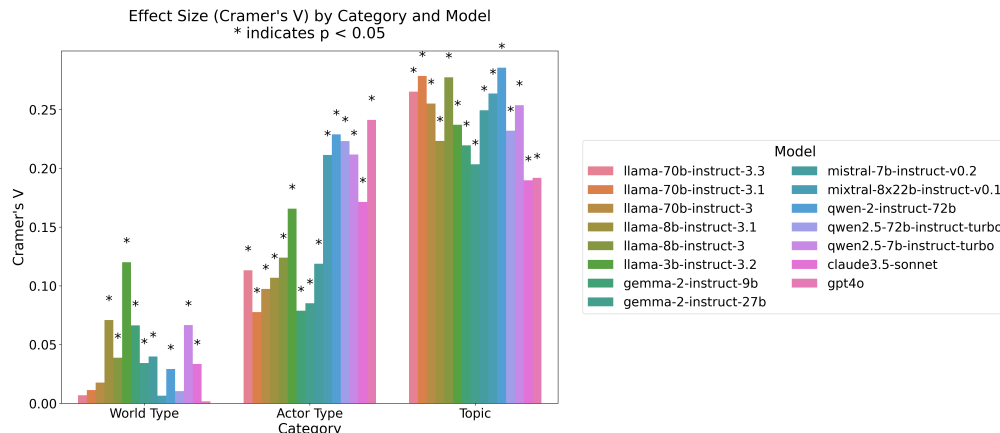


Figure 8: Effect size per model and per topic as measured by Cramer's V. Across model families and sizes, Topic consistently has the largest effect on decision, with the only outlier being GPT4o for which Actor Type has the largest impact.

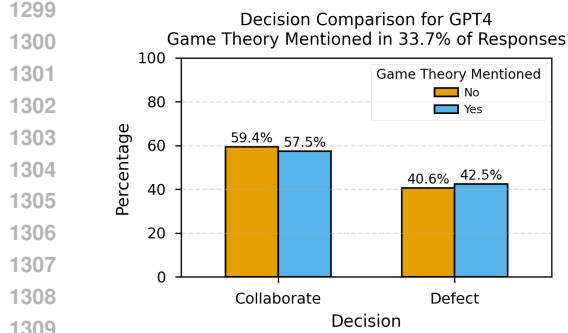
F GAME RECOGNITION

In response to each scenario, Claude, Llama, and GPT-4 were required to provide both a decision and its justification. The models frequently incorporated game theory concepts in their reasoning, though the frequency varied considerably by model (ranging from 16.5% to 70.5%). These percentages were calculated by using Claude 3.5 Sonnet to analyze the reasoning chains. For each justification, Claude was prompted with the following question:

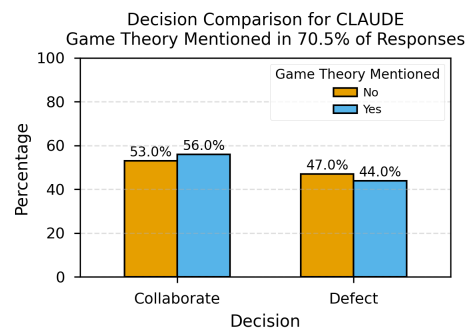
PROMPT: DOES THIS TEXT EXPLICITLY MENTION THE PRISONER'S DILEMMA OR GAME THEORY? RESPOND ONLY WITH <YES> OR <NO> FOLLOWED BY THE RELEVANT SENTENCE(S). HERE IS THE TEXT: TEXT

The models' tendency to reference game theory suggest they recognize these scenarios as formal games and use this framework in their decision-making process. In Appendix F, we show how

1296 the cooperate/defect proportion changes depending on whether game theory was mentioned in the
 1297 justification. We see small, yet non-trivial changes to the proportion of cooperation.
 1298



1310 (a) GPT-4o



1331 (b) Claude

1332 G HEATMAPS

1333

1334

1335

1336

1337

1338

1339

1340

1341

1342

1343

1344

1345

1346

1347

1348

1349

1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403

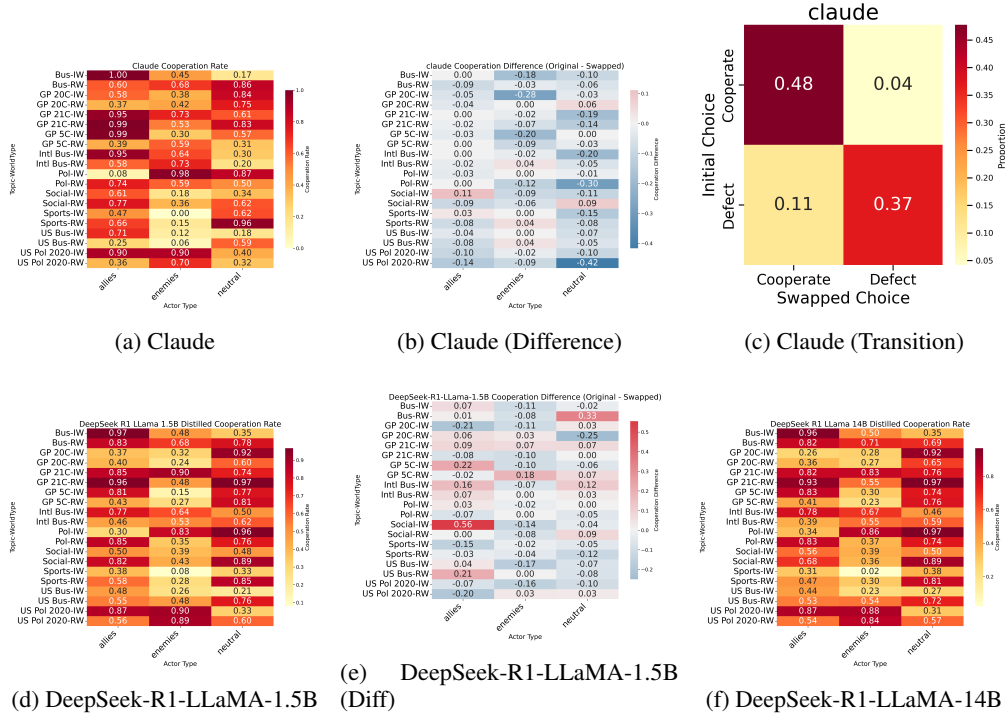


Figure 10: Distribution of decisions: Claude and DeepSeek-R1 models.

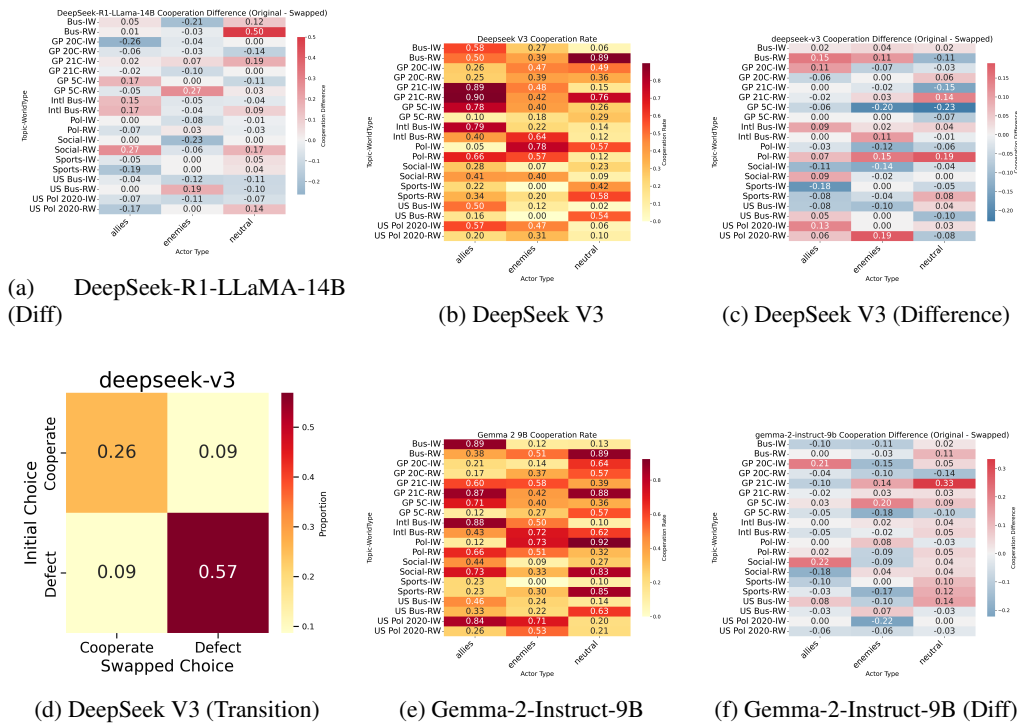


Figure 11: Distribution of decisions: DeepSeek and Gemma-2-9B models.

1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457

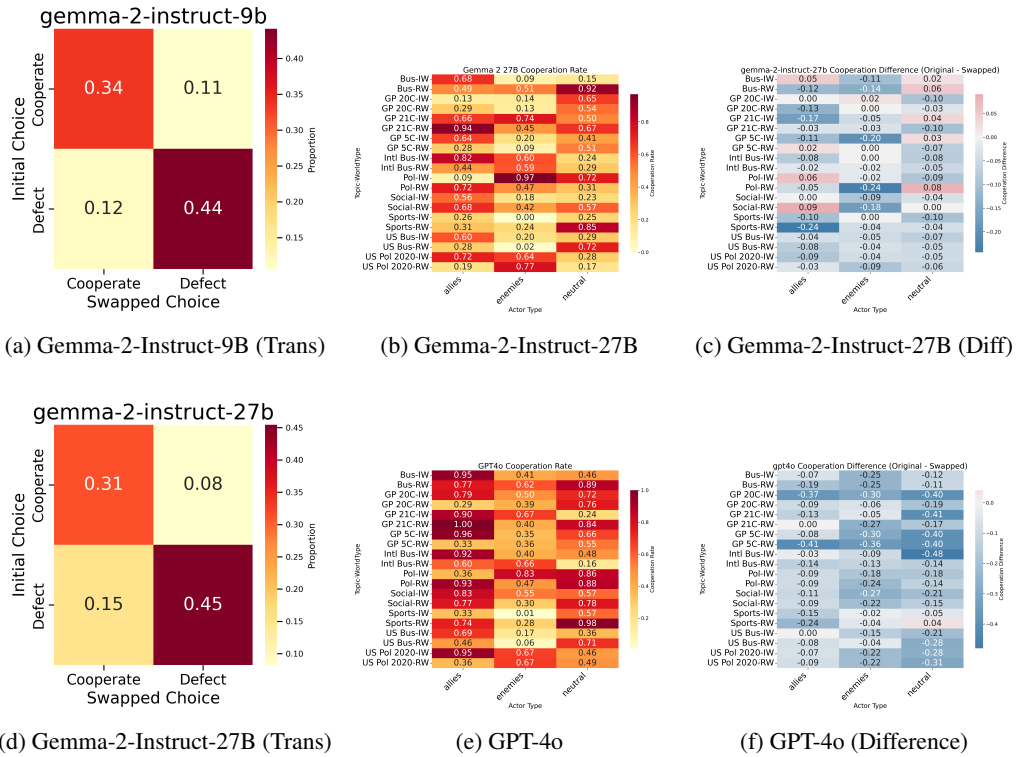


Figure 12: Distribution of decisions: Gemma-2 and GPT-4o models.

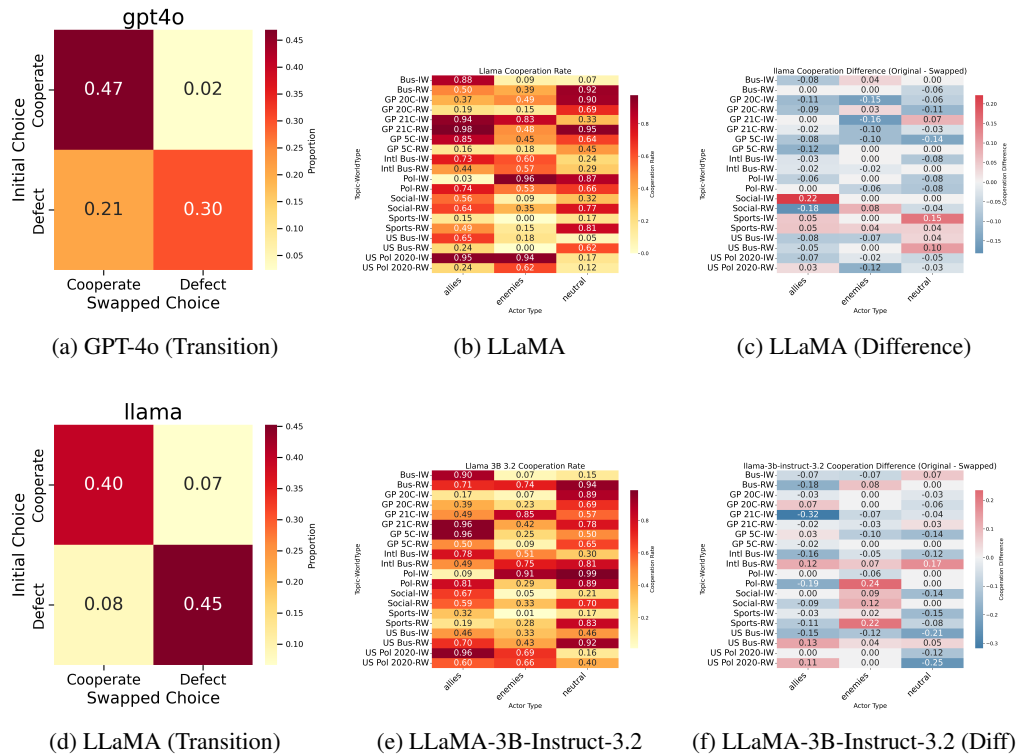


Figure 13: Distribution of decisions: GPT-4o and LLaMA base models.

1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511

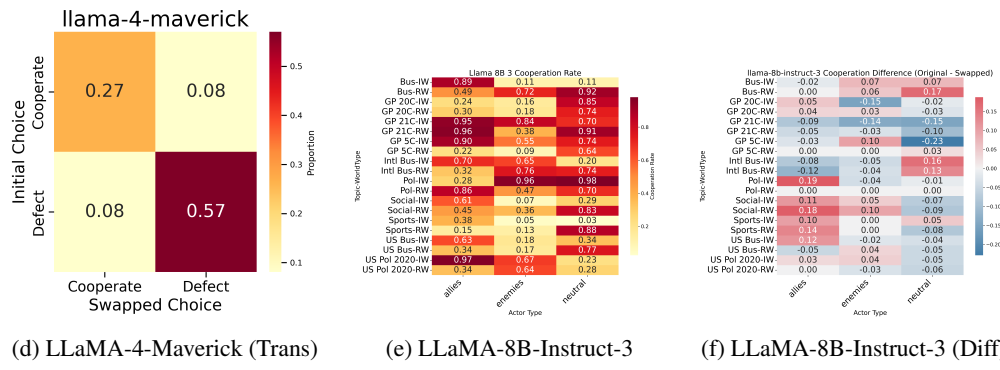
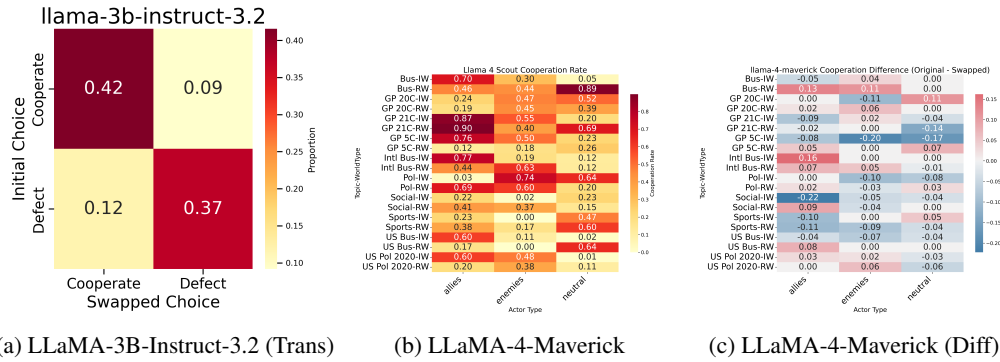


Figure 14: Distribution of decisions: LLaMA-3B, LLaMA-4-Maverick, and LLaMA-8B-3 models.

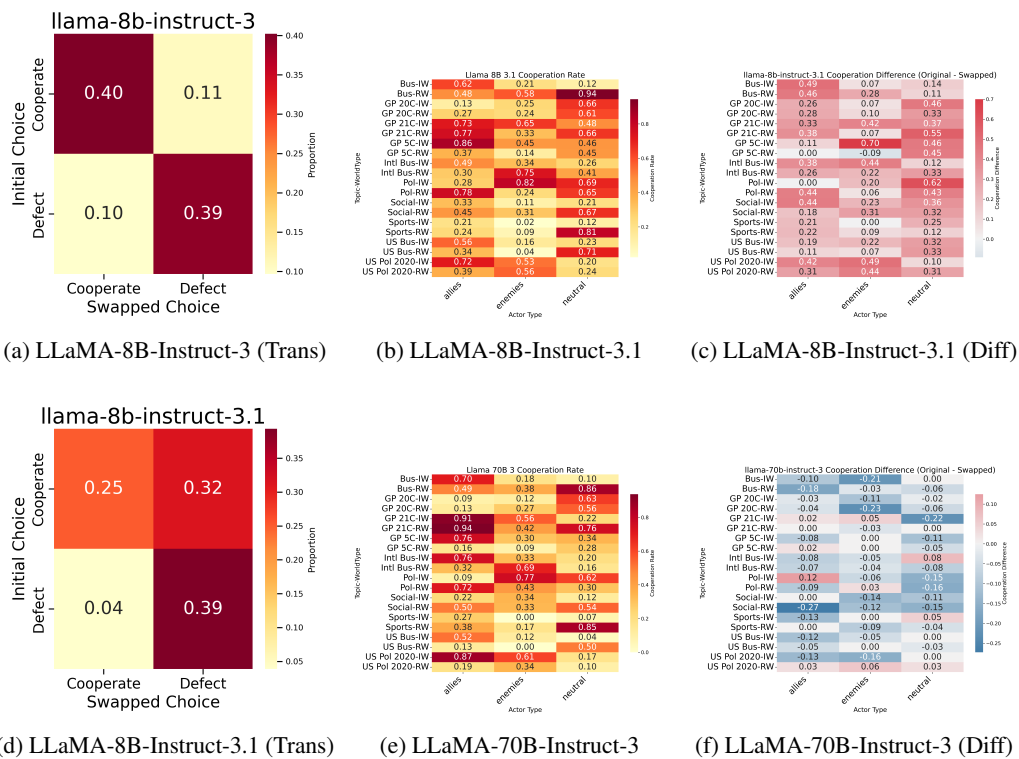
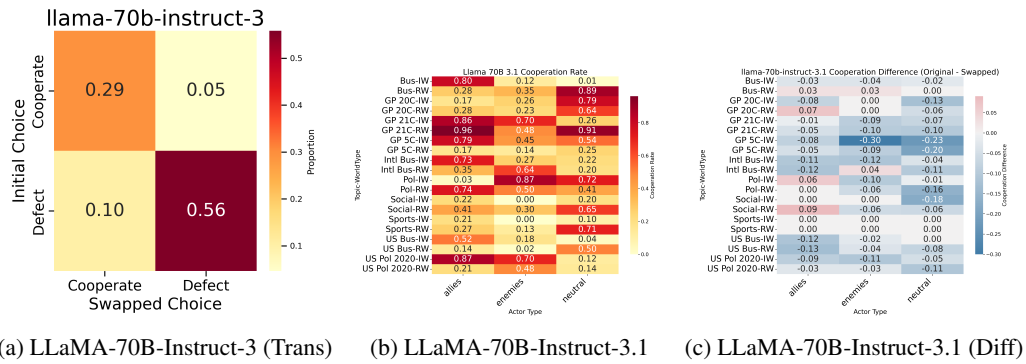


Figure 15: Distribution of decisions: LLaMA-8B-3.1 and LLaMA-70B-3 models.

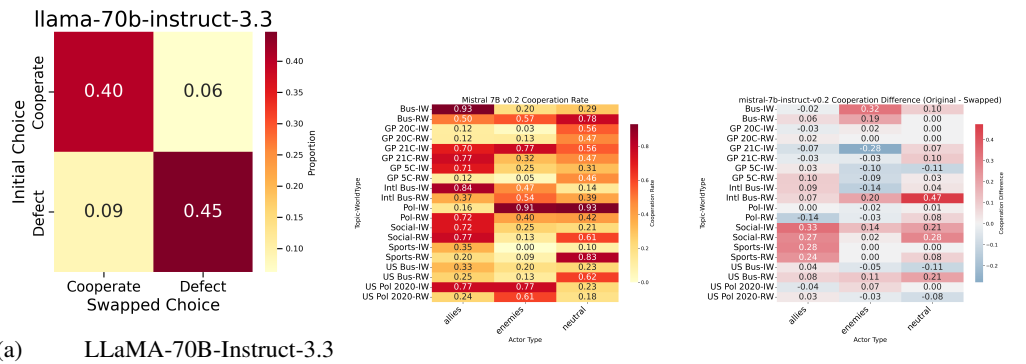
1512
1513
1514
1515
1516
1517
1518
1519
1520
1521
1522
1523
1524
1525
1526
1527
1528
1529
1530
1531
1532
1533
1534
1535
1536
1537
1538
1539
1540
1541
1542
1543
1544
1545
1546
1547
1548
1549
1550
1551
1552
1553
1554
1555
1556
1557
1558
1559
1560
1561
1562
1563
1564
1565



(a) LLaMA-70B-Instruct-3 (Trans) (b) LLaMA-70B-Instruct-3.1 (c) LLaMA-70B-Instruct-3.1 (Diff)

(d) LLaMA-70B-Instruct-3.1 (Trans) (e) LLaMA-70B-Instruct-3.3 (f) LLaMA-70B-Instruct-3.3 (Diff)

Figure 16: Distribution of decisions: LLaMA-70B-Instruct models (3.1 and 3.3).



(a) LLaMA-70B-Instruct-3.3 (Trans) (b) Mistral-7B-Instruct-v0.2 (c) Mistral-7B-Instruct-v0.2 (Diff)

(d) Mistral-7B-Instruct-v0.2 (Trans) (e) Mistral-7B-Instruct-v0.3 (f) Mistral-7B-Instruct-v0.3 (Diff)

Figure 17: Distribution of decisions: LLaMA-70B-3.3 and Mistral-7B models.

1566
1567
1568
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578
1579
1580
1581
1582
1583
1584
1585
1586
1587
1588
1589
1590
1591
1592
1593
1594
1595
1596
1597
1598
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619

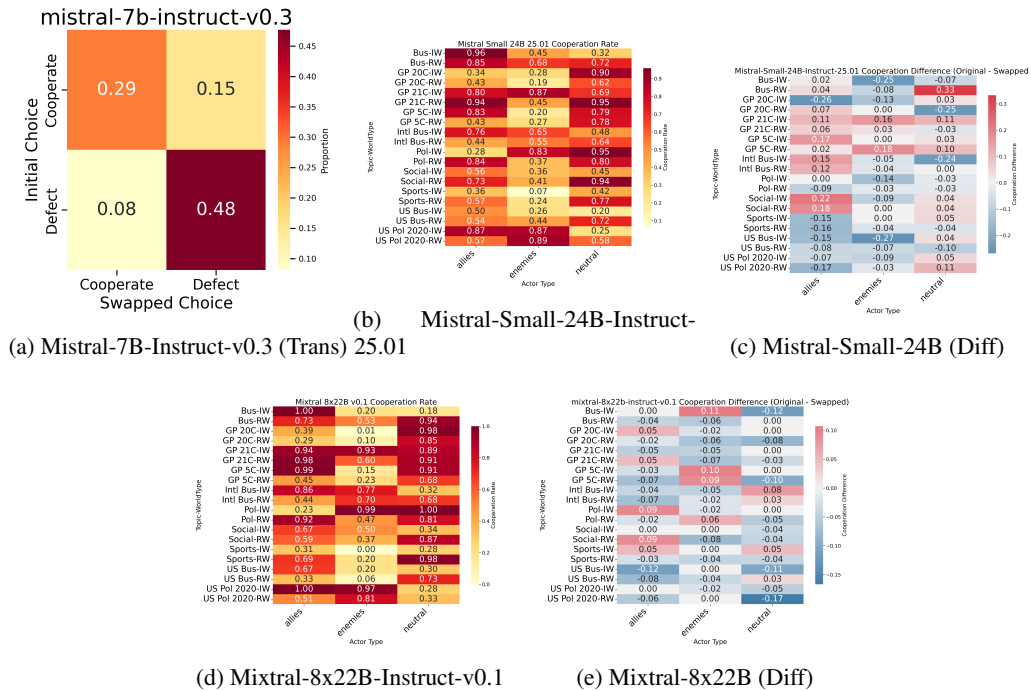


Figure 18: Distribution of decisions: Mistral-Small-24B and Mixtral-8x22B models.