# Large language models replicate and predict human cooperation across experiments in game theory

**Andrea Cera Palatsi**[1], **Samuel Martin-Gutierrez**[2], **Ana S. Cardenal**[3], **and Max Pellert**[1,*]

[1]Department for Computational Social Sciences and Humanities, Barcelona Supercomputing Center
[2]Grupo de Sistemas Complejos, Universidad Politécnica de Madrid
[3]School of Law and Political Science, Universitat Oberta de Catalunya
[*]max.pellert@bsc.es

## ABSTRACT

Large language models (LLMs) are increasingly used both to make decisions in domains such as health, education and law, and to simulate human behavior. Yet how closely LLMs mirror actual human decision-making remains poorly understood. This gap is critical: misalignment could produce harmful outcomes in practical applications, while failure to replicate human behavior renders LLMs ineffective for social simulations. Here, we address this gap by developing a digital twin of game-theoretic experiments and introducing a systematic prompting and probing framework for machine-behavioral evaluation. Testing three open-source models (Llama, Mistral and Qwen), we find that Llama reproduces human cooperation patterns with high fidelity, capturing human deviations from rational choice theory, while Qwen aligns closely with Nash equilibrium predictions. Notably, we achieved population-level behavioral replication without persona-based prompting, simplifying the simulation process. Extending beyond the original human-tested games, we generate and preregister testable hypotheses for novel game configurations outside the original parameter grid. Our findings demonstrate that appropriately calibrated LLMs can replicate aggregate human behavioral patterns and enable systematic exploration of unexplored experimental spaces, offering a complementary approach to traditional research in the social and behavioral sciences that generates new empirical predictions about human social decision-making.

The capabilities of large language models (LLMs) enable their use as digital twins to be both studied and used for inference in the social and behavioral sciences. LLMs as computational complements for humans offer the promise of first replicating experiments originally conducted with humans and then extending them to previously unexplored areas, making the experimental process much faster and cheaper while operating under less stringent ethical constraints. The algorithmic fidelity[1] of such synthetic approaches hinges on the ability to robustly elicit valid behaviors from LLMs. While their training data potentially contains vast implicit information on human patterns of experience and behavior, researchers face many modeling choices and challenges in their attempts to reliably extract those patterns. To simulate social experiments, we first need to systematically assess whether LLMs can replicate empirically observed human patterns to a satisfactory degree.

By comparing the behavior of LLMs subjected to the same conditions as in experiments originally conducted with human participants, we gain crucial insights into machine behavior[2]. This line of inquiry is becoming increasingly relevant as LLMs are deployed as autonomous agents in domains such as medicine, law, and education[3–9]. When such agents make decisions, it is crucial to understand the factors influencing their choices. Certain LLMs may exhibit systematic differences in decision-making compared to humans, displaying biases relative to human norms. Even within the same model, variations such as different prompting strategies can lead to undesired outcomes.

In addition, we aim to further our understanding of general patterns observable across models. Universal features uncovered in LLMs could in turn help us learn more about human behavior. The realism of initial attempts to use LLMs as "cognitive models"[10] remains debated[11] and part of the model performance can be due to Clever Hans effects[12]. Rather than attempting to simulate individual participants, which presents substantial challenges, we focus on the more tractable goal of capturing empirical *macro-level* behavioral patterns, an approach that, while not without its own challenges, provides an effective means for understanding collective human behavior. Validation by successfully replicating existing human experiments is only the first step: We want to build on rich empirical evidence to ensure our simulations match it, while also being able to derive testable hypotheses about experiments that have *not yet been conducted*.

The methodological freedom that comes with modeling using LLMs also introduces complexity, making it essential to work within well-defined domains to evaluate approaches clearly. Classical games from game theory are well-suited for this purpose, as they offer 1) analytical (often "optimally rational") solutions, 2) ample empirical data of humans playing such games and 3) mathematically well-defined, simple tasks with clear instructions that can be readily translated to natural language. Furthermore,

despite their simplicity, these game-theoretic scenarios are ubiquitous and offer powerful analytic tools for analyzing many different complex phenomena in our social worlds relevant to contexts where LLM agents may operate.

## Examples of LLMs playing games

Existing research has examined LLMs in various game theoretical settings including the Dictator, Ultimatum, Trust Game and the Prisoner's dilemma as well as classic behavioral experiments such as Wisdom of Crowds and the Milgram obedience experiment (for details see A.1).

The earliest work in this domain tested LLMs (different models from the GPT series) with the Dictator Game as well as with other decision-making scenarios[13]. When given no persona instructions, the most advanced model tested (GPT-3 text-davinci-003) consistently selected efficiency-maximizing outcomes and did not vary its responses to the same game when asked repeatedly. The authors instructed LLMs in different conditions to be inequity-averse (care about fairness between players), efficient (care about the total payoff of both players) and selfish (care about their own payoff). The models responded by following these instructions precisely when capable, though less advanced GPT-3 models could not modulate their behavior based on instructions. When the model (GPT-3 text-davinci-003) received no instructions about how to value fairness, efficiency, or self-interest, it consistently chose efficient outcomes that maximized total payoff. However, to replicate actual human behavior patterns, a synthetic population needed to be composed of 52% selfish agents, 32% efficiency-focused agents, and only 15% inequity-averse agents, suggesting that the model's default preferences, inclined towards efficiency and increasing common benefit, differ from human tendencies, mostly associated with selfish behaviors.

The same authors also simulated a scenario in which a store increased the price of snow shovels after a snowstorm and found that LLMs were more likely to reject larger price increases. Prompting models to have different political ideologies changed their responses (e.g. left-wing personas consistently rated price increases as unfair, while moderate and libertarian personas found smaller increases acceptable). Another scenario had models choosing between different options with one that was presented as the status quo. Models, similar to humans, preferred the status quo. A final scenario had models acting as employers choosing between candidates with different experience levels. When a minimum wage was imposed, LLMs shifted towards hiring more experienced applicants, a pattern also found among humans.

A similar line of research tested eight GPT model variants, from text-ada-001 to GPT-4[14]. In the Ultimatum Game, text-davinci-002 did not accept unfair proposals below 10% of the endowment, and showed high acceptance rates for offers of 50% or more, a gradual pattern that appears in humans, showing behavior consistent with human tendencies to reject unfair proposals. LLMs also played the Milgram obedience experiment, with researchers creating an alternative "sleepiness and driving" scenario to prevent the models from replicating memorized descriptions of the original experiment. In the end, text-davinci-002 displayed the same tendency of obedience that was originally observed in humans, though at slightly higher rates (75% vs. 65% in the original human study). While text-davinci-002 successfully replicated human behavioral patterns in these experiments, more recent models like GPT-4 were only tested in the Wisdom of Crowds study, where they exhibited "hyper-accuracy distortion" by providing inhumanly accurate answers (often perfectly correct) rather than the varied, imperfect estimates typical of humans. This prevented the Wisdom of Crowds aggregation effect from functioning properly, revealing how alignment procedures that improve factual accuracy can create systematic differences from human cognition.

Mei et al.[15] had LLMs play games and administered the Big 5 personality questionnaire to the models themselves, having them answer psychological survey questions directly. Apart from measuring behavioral traits, they performed Turing experiments, selecting individual responses to games and comparing them with human responses. For the Big 5 personality test, they compared the answers of two models with human distributions. They confirmed that responses of the same model did not present much variation, similar to the consistency a single human subject would show if repeatedly queried. The Big 5 results of the advanced model were the most similar to humans, becoming statistically indistinguishable. For the games, GPT-4 passed the behavioral Turing test in most games, but showed detectably different patterns in the Trust Game, where it demonstrated greater trust and generosity than typical humans, consistently avoiding purely selfish strategies (like giving $0) that humans frequently chose. Humans exhibited more extreme and varied actions across the behavioral spectrum. In the Prisoner's Dilemma, it was more cooperative than the median human, who showed greater inclination to defect. The advanced model showed adaptive behavior, modifying its responses based on previous experience "as if" learning from interactions. In the Prisoner's Dilemma, GPT-4 modified its responses depending on the past, showing conditional cooperation strategies similar to those observed in humans.

In summary, several stylized facts emerge from existing research on LLMs playing classical game-theoretic scenarios. Advanced models can exhibit behavioral patterns similar to humans in experimental settings, though when differences arise, LLMs typically demonstrate greater cooperation than human participants. A key methodological challenge is that each model instance effectively acts as a single subject, requiring deliberate strategies to recreate the diversity observed in human populations. While extended interviews show the greatest potential for endowing LLMs with distinct personalities, simpler approaches using demographic attributes or brief descriptions can also generate meaningful variation. Notably, LLMs tend to

attribute higher knowledge and accuracy to others than is empirically typical (a form of hyper-accuracy bias), and advanced models often exhibit payoff-maximizing strategies that benefit both players. Additionally, LLMs demonstrate learning capacity and can adapt their strategies based on interaction history. Finally, most existing studies rely on commercial models accessed via closed, black-box APIs, limiting transparency and reproducibility.

We extend these earlier approaches in several ways. First, we conceptually integrate our work into the larger research program of machine psychology[16,17]. Second, by working within the formal framework of game theory, we can derive analytical solutions (such as Nash equilibria) that serve as benchmarks for comparison, allowing us to systematically evaluate LLM behavior against both human empirical patterns and theoretical predictions. Third, rather than solely replicating known empirical patterns from existing human experiments, we use LLMs to simulate novel experimental conditions, generating empirically testable hypotheses for future human studies. Fourth, we employ open-source models to ensure full transparency and reproducibility of our findings.

## Choice of models

We focus on widely used, open-source models that are both archetypal of current LLM capabilities and likely to be deployed as agents in many different contexts. We selected three medium-sized models that balance performance with computational feasibility while ensuring architectural diversity. We derive our main results using **Llama-3.1-8B-Instruct** by Meta AI[18], as justified in Section 1. To assess robustness across model architectures, we repeat the exact same analyses with **Mistral-7B-Instruct-v0.3** by Mistral AI[19] and **Qwen2.5-7B-Instruct** by Alibaba Cloud[20,21]. Additionally, we employ Qwen2.5-7B-Instruct for the separate tasks of answer extraction and logical verification (for details see Subsection 3.4 and Table 3).

For simplicity, we refer to these models as Llama, Mistral, and Qwen throughout the remainder of this paper.

## Digital twinning of game theoretical experiments

We based our LLM setup on the experiments of Poncela-Casasnovas et al.[22], who collected data from more than five hundred participants playing dyadic games. Participants were recruited from the general audience of the game festival DAU Barcelona (December, 2014) to ensure a diverse sample, resulting in an average age of 31.3 (SD = 14.3) and 64.5% females. The authors created an application on tablets that allowed participants to play games with random opponents and payoff configurations randomly selected from a pre-set parameter grid. The experiment was conducted across multiple sessions over a period of two days. After each game, experimenters paid participants with lottery tickets according to the points they scored. Each game consisted of a single round in which both players made one decision simultaneously, without receiving feedback about the opponent's choice until after both had decided.

Each player could choose to cooperate or defect (encoded as colors to prevent wording bias for the human participants). Depending on the choices of both players, there are four possible outcomes (Table 1). If a player chooses C, they receive payoff $R$ if the opponent also chooses C, or payoff $S$ if the opponent chooses D. If a player chooses D, they receive payoff $T$ if the opponent chooses C, or payoff $P$ if the opponent also chooses D.

**Table 1. Payoff matrix for strategies C (Cooperate) and D (Defect).** Each cell shows the payoff as (Player 1, Player 2). The parameters represent: $R$: **r**eward for mutual cooperation; $P$: **p**unishment for mutual defection; $T$: **t**emptation to defect; $S$: **s**ucker's payoff.

|  |  | Player 2 | |
|---|---|---|---|
|  |  | C | D |
| **Player 1** | C | $(R,R)$ | $(S,T)$ |
|  | D | $(T,S)$ | $(P,P)$ |

In the experiment, the value of $R$ was kept constant at 10 points and the value of $P$ at 5 points, ensuring that mutual cooperation always yields a higher payoff than mutual defection. The parameter $S$ ranged from 0 to 10, while $T$ ranged from 5 to 15, with both taking integer values. This created 11 possible values each parameter, leading to 121 distinct games.

These parameter combinations constitute a generalized framework that encompasses four classical game types:

- HG (Harmony Game): $S > P$, $R > T$, or equivalently $S > 5$ and $T < 10$. Cooperation is the dominant strategy.

- SG (Snowdrift Game): $T > R > S > P$, or equivalently $T > 10 > S > 5$. The best individual outcome is to defect when the opponent cooperates, but the worst outcome is mutual defection.

- SH (Stag Hunt Game): $R > T \geq P > S$, or equivalently $10 > T \geq 5 > S$. The best outcome is mutual cooperation, but the worst outcome is to cooperate when the opponent defects.

- PD (Prisoner's Dilemma): $T > R > P > S$, or equivalently to $T > 10 > 5 > S$. The best individual outcome is to defect when the opponent cooperates, and the worst is to cooperate when the opponent defects.

To analyze human participants' behavior across games with uniformly randomly selected parameters, the authors used K-means clustering to identify 5 distinct behavioral phenotypes that characterize player strategies. The distribution of humans assigned to each phenotype was: 20% Optimist, 21% Pessimist, 30% Envious, 17% Trustful, and 12% Undefined. Notably, all these phenotypes deviate from the Nash equilibrium and therefore fall outside the bounds of classical rationality. The authors identified a decision rule for each phenotype:

- **Optimist.** Maximizes the best possible outcome. Cooperates if $R > T$, assuming the opponent will cooperate.

- **Pessimist.** Maximizes the worst possible outcome (minimax strategy). Cooperates if $S > P$.

- **Envious.** Seeks to maximize relative advantage over the opponent. Cooperates if $S \geq T$.

- **Trustful.** Always cooperates, representing the strongest inclination toward maximizing joint payoffs.

- **Undefined.** Exhibits no consistent strategy, cooperating and defecting randomly.

To validate these decision rules, the authors calculated, for each phenotype and game, the difference between predicted cooperation rates (based on the theoretical rules) and observed cooperation rates (from the empirical data). These differences were then normalized by the standard deviation of the human sample and averaged across all games, yielding a mean deviation for each phenotype. Since none of the phenotypes exhibited an average deviation beyond the 99% confidence interval (2.575 SD units), the authors concluded that the theoretical cooperation patterns were statistically consistent with the empirical human data.

We adopted and adapted this experimental framework for use with LLMs, employing multiple games rather than a single game to identify cross-game behavioral patterns and characterize behavioral differences. Specifically, we simulated the same games played in the original experiment, enabling direct comparison between human and LLM behavior. However, we went beyond mere replication by introducing novel games designed to explore behavioral patterns that have not yet been tested with humans. In doing so, we employ LLMs as a predictive proxy for human behavior, demonstrating their potential to forecast how humans might behave in unexplored strategic settings, which represents a key application of AI for scientific discovery.

Our work builds upon and extends earlier attempts to subject LLMs to classical game-theoretic scenarios. While these earlier studies are valuable as proof-of-concept demonstrations, they typically lack thorough calibration against empirical human data. We advance this research program in several ways. First, we systematically validate our approach using multiple open-source models rather than relying on a single type of commercial model, enabling assessment of robustness across different architectures. Second, we develop and validate a rigorous prompting methodology informed by empirical testing, incorporating logical verifiers and other quality control mechanisms to ensure response validity. Third, we use this validated setup not merely to replicate known patterns, but to systematically expand the experimental parameter space, generating empirically testable hypotheses for novel games that can guide future human experiments.

# 1 Results

We present the results visually as cooperation matrices. In these matrices, each cell represents the average cooperation rate for a specific game defined by a particular combination of $S$ and $T$ values.

## 1.1 Answer extraction refinement strategies

We developed four progressively more complex methods for extracting answers from the models (for details see Subsection 3.4):

- **Simple extraction.** After providing the instructions, we directly ask the tested LLM to state its chosen strategy by writing only the corresponding label (e.g., *A* for defection and *B* for cooperation).

- **Double extraction.** We ask the tested LLM for a long answer, then use Qwen to extract the label of the chosen strategy from this longer response.

- **Multi-step extraction.** We ask the tested LLM for a long answer using a prompt that guides the model through logical reasoning steps, then use Qwen to extract the strategy label.

- **Logical verifier extraction.** We ask the tested LLM for a long answer using a prompt that guides the model through logical reasoning steps (for details see Subsubsection 3.4.4). Before extraction, we apply Qwen as a logical verifier to filter the responses, and then use Qwen again to extract the strategy label.

The left-hand side of Figure 1 demonstrates how each step of our refinement procedure progressively reduces noise in the cooperation patterns of the Llama model, revealing a pattern that closely resembles the empirical human cooperation matrices: First, when using only "Simple Extraction", cooperation appears almost random. Second, upon adding "Double Extraction", we start observing increased cooperation in the Harmony Game ($S > 5$, $T < 10$). Third, with "Multi-step Extraction", two distinct regions emerge: games where $S \geq T$ exhibit increased cooperation, while games where $T > R$ exhibit decreased cooperation. With the final layer of complexity ("Logical Verifier"), we closely replicate the original empirical patterns of humans shown in Panel A of Figure 2.
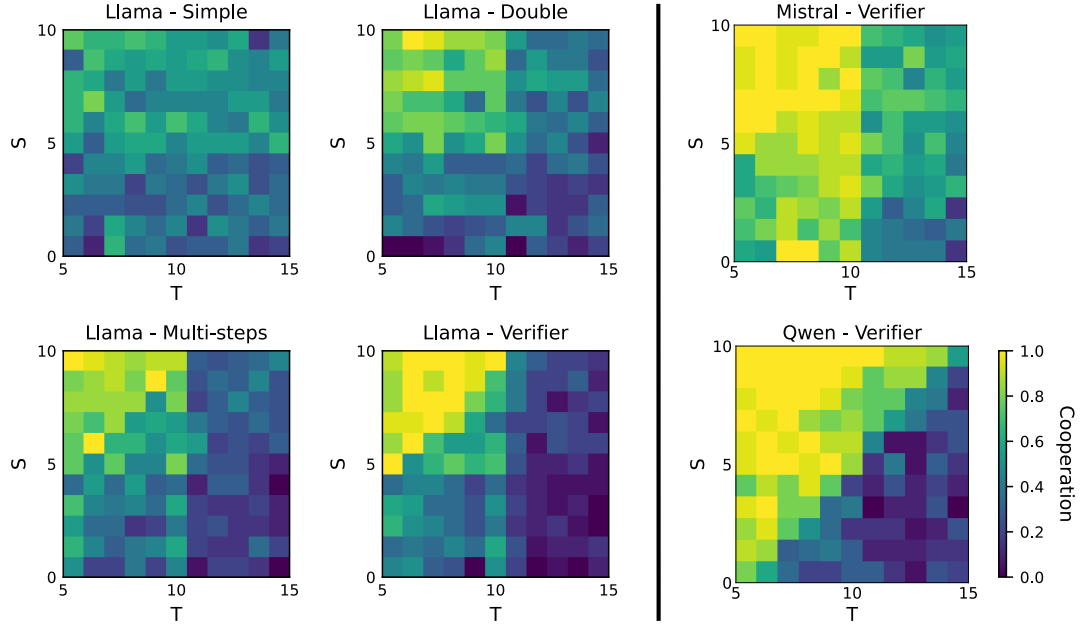


**Figure 1.** **Effect of different extraction strategies on Llama's cooperation matrix and comparison with other models.** Left panel: Progressive improvements in revealing human empirical patterns in Llama's cooperation matrix as extraction complexity increases. (For the other two models, see Subsection A.1.) In the final step (bottom right), Llama shows an average cooperation of 0.402. Right panel: Average cooperation is 0.722 for Mistral and 0.568 for Qwen. Color scales indicate average cooperation from 0 (purple: no cooperation) to 1 (yellow: full cooperation).

## 1.2 Differences of cooperation across LLMs

Figure 1 shows the cooperation patterns of all three models (Qwen and Mistral in the right panel and Llama in the bottom right of the left panel). Each model exhibits distinct behavioral characteristics: Llama is most similar to the patterns found in the empirical data of human participants, showing the same qualitative features in its cooperation matrix. Mistral displays a vertical line that separates the regions where $T > R$ (less cooperative), from regions where $T \leq R$ (more cooperative). This pattern suggests an eagerness towards earning the highest possible reward, which resembles the human "optimist" phenotype. Qwen closely follows the Nash equilibrium (Panel C, the non-transparent inset of Figure 2).

To quantify these differences, Table 2 presents the mean squared displacement (MSD) and Pearson's $r$ comparing the three models with both human data and the Nash equilibrium. The metrics confirm our qualitative observations: Llama shows the highest similarity to humans (MSD = 0.031, $r = 0.89$), Mistral is similar to neither humans nor the Nash equilibrium and Qwen exhibits the highest similarity to the Nash equilibrium (MSD = 0.036, $r = 0.93$).

As Llama most closely aligns with human behavioral patterns, we use it in subsequent analyses to replicate and extend the original experiments with humans. For further discussion of the other two other models and robustness across model architectures, see Subsection A.3 in the Appendix.

|  | Human | | Nash | |
| --- | --- | --- | --- | --- |
|  | **MSD** | *r* | **MSD** | *r* |
| **Llama** | **0.031** | **0.89** | 0.089 | 0.77 |
| **Mistral** | 0.091 | 0.70 | 0.182 | 0.60 |
| **Qwen** | 0.065 | 0.79 | **0.036** | **0.93** |

**Table 2. Metrics for model comparisons.** Mean squared displacement (lower is better) and Pearson's *r* (higher is better) comparing the cooperation matrices of the three LLMs with those of human participants and the Nash equilibrium. Llama is closest to the human data, while Qwen exhibits the highest similarity to the Nash equilibrium.

### 1.3 Replicating the original games

For this analysis, we focus on the non-transparent regions of Figure 2 (delineated with black borders), which show only the parameter combinations played in the original experiment.

Panel A (upper part) of Figure 2, shows the human behavioral patterns from the original experiment. In the Harmony Game (HG) region there is a clear area defined by $S \geq T$ that displays distinctly high cooperation. The remainder of the matrix shows similar cooperation levels, with slightly lower rates in the Prisoner's Dilemma (PD) region. Llama's behavior (Panel B of Figure 2) shows qualitatively similar cooperation patterns compared to humans exhibiting high cooperation in the region where $S \geq T$. In contrast, the region where $T > R$ displays lower cooperation levels. The statistical measurements in Table 2 quantify the similarity between Llama and the human cooperation matrix. The Mean Squared Displacement (MSD) of 0.031 is small considering the possible maximum of 1, while Pearson's *r* of 0.89 indicates a strong linear correlation between the two matrices.

The Nash equilibrium cooperation matrix (Panel C of Figure 2) represents rational behavior. The rational choice for HG is full cooperation. For the Snowdrift Game (SG), we find a mixed equilibrium, where both strategies are expected in different proportions. For the Stag Hunt (SH) game, there is a boundary that separates two sub-regions depending on the stable strategy, while for the Prisoner's Dilemma (PD) the rational choice is to defect. Among the three cooperation matrices for the original game parameters shown in Figure 2, the Nash Equilibrium matrix displays the highest average cooperation and Llama the lowest. Therefore, on average, Llama cooperates less than humans, who in turn cooperate less than theoretically expected.

### 1.4 Extending the experimental grid: novel in-silico games

Using the exact same setup as before, we simulated Llama playing games with parameter combinations outside the regions covered in the original experiments with human participants. Figure 2 shows this full cooperation matrix with novel combinations of $S$ and $T$ (shown as transparent regions) alongside simulated games already played in the experiments (non-transparent, delineated with a black border).

Llama exhibits high cooperation when $S \geq T$, extending the diagonal $S = T$ pattern into the enlarged cooperation matrix and demonstrating the importance of this parameter condition. Within the region where $S \geq T$, there is no apparent effect of the values $P$ and $R$ on the tendency to cooperate. However, when $S$ and $T$ have values near 0, the separation becomes less distinct, displaying a random-like cooperation pattern when $S < 5$ and $T < 5$. In the other half of the matrix, where $S < T$, the value of $R$ exerts a clear influence: there is a marked difference between regions where $T > R$ (exhibiting clearly lower cooperation) and $T \leq R$. Examining horizontal cross-sections at constant $S$ reveals that increasing $T$ reduces cooperation, though not linearly. Similarly, vertical cross-sections at constant $T$ show that increasing $S$ increases cooperation.

If we were instead to use the Nash equilibrium solution to predict the novel games (Panel C of Figure 2), we would find two pure regions of total defection and total cooperation separated by a boundary where the average cooperation is 0.5, along with a third region is characterized by mixed equilibrium. Both the Llama simulations and the Nash equilibrium predictions provide empirically testable hypotheses about human behavior in unexplored parameter regions. To decisively evaluate which framework—classical rationality or LLM-based modeling—better anticipates real human behavior, we have preregistered[1] the precise experimental setup that will allow future experiments to establish which approach shows more similarity with actual human behavior in these novel games.

## 2 Discussion

In this study, we created a digital twin of a classic game-theoretic experiment using large language models (LLMs). By systematically replicating the original setup across several models (Llama, Qwen, and Mistral) we evaluated their ability to reproduce empirically observed human behaviors. Llama closely mirrors human cooperation patterns, and Qwen aligns with
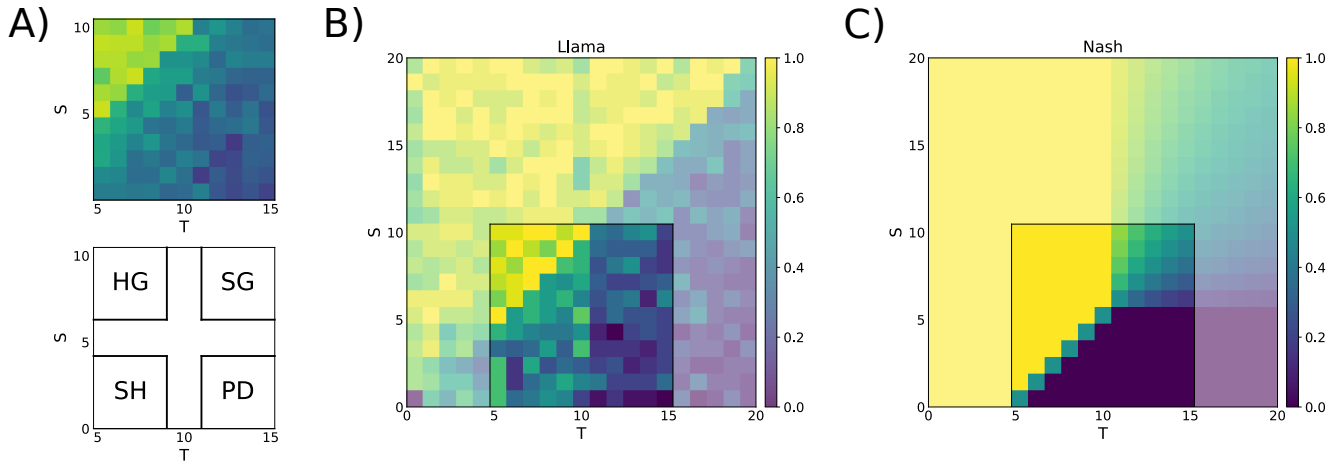
---

[1] https://aspredicted.org/fe6z2k.pdf

**Figure 2.** **Cooperation matrices of human participants, Llama and the Nash equilibrium, extended by simulations.**
Each tile shows a game with parameter combination of $(S, T)$, with average cooperation ranging from 0 (purple: no cooperation) to 1 (yellow: full cooperation). Panel A (top): Human cooperation matrix from the original experiments[22] with an average cooperation of 0.480. Panel A (bottom) shows that each quadrant of the parameter grid corresponds to a standard game in game theory: HG (Harmony Game), SG (Snowdrift Game), SH (Stag Hunt Game), and PD (Prisoner's Dilemma). Games outside the quadrants are non-standard. Panel B: Cooperation matrix Llama simulations. The non-transparent region delineated with a black border shows simulations with parameter combinations that were already present in the original human experiments (average cooperation of Llama of 0.402 in this region). The transparent region shows combinations of $S$ and $T$ that were *not yet played* by humans. Panel C: Matrix showing the rational choices predicted by Nash equilibrium (see Subsection 3.6) covering both original and novel games. Average cooperation in the original games (non-transparent region) is 0.500. The simulations of the extended parameter grid generate testable hypotheses that can be validated through preregistered experiments to establish which extrapolation method for novel games agrees better with empirical human behavioral patterns.

the theoretical Nash equilibrium. Extending the parameter space beyond the range explored with human participants, we further derived predictions that can be directly tested in future behavioral experiments.

Behavioral patterns observed in Llama reveal several interesting features when replicating the original games. Notably, the region where $T > R$ displays lower cooperation levels, potentially reflecting the model's sensitivity to the temptation payoff exceeding the mutual cooperation reward. This pattern suggests that, like humans, Llama may be responding to the relative attractiveness of defecting when the opponent cooperates, even though this does not align with optimal collective outcomes.

Our extension of the parameter space, while modest compared to all logically possible combinations, was strategically designed to systematically explore regions adjacent to the original experimental grid and to make it symmetric. This approach allows us to test whether patterns observed in the original games generalize smoothly to nearby parameter values or whether discontinuities emerge at certain boundaries. The original games represent only a small subset of the possible experimental space, and our extensions provide a principled way to extrapolate beyond known regions while maintaining interpretability.

A key distinction emerges between prediction using the analytical Nash equilibrium solution versus the LLM-based approach. The analytical solution predicts regions with intermediate cooperation rates (around 0.5) corresponding to mixed equilibria, particularly visible in the Snowdrift Game region of Panel C in Figure 2. These smooth gradients do not appear in Llama's behavior (Panel B), which instead exhibits overall more discrete, binary-like cooperation patterns. Mixed equilibria would require players to randomize their strategies in precise proportions, which is a sophisticated form of strategic reasoning that can be considered psychologically implausible for most players. Instead, both humans and LLMs likely tend to employ simpler decision heuristics based on comparing payoffs or identifying dominant features of the game structure. The closer alignment of Llama with human patterns rather than Nash predictions suggests that LLMs may capture the psychological and heuristic processes underlying human decision-making in ways that purely analytical Nash equilibrium solutions cannot. These findings are consistent with the theory of bounded rationality, which emphasizes that models of social behavior must account for the limited access to information and, crucially, the computational capacities that are actually possessed by people[23]. In cognitive science, a *simplicity principle* highlights our tendency to adopt simple strategies when faced with complex tasks[24]. Empirically, researchers have found that we follow simple heuristics when choosing social connections[25] and that limited cognitive capacities explain the structure of our social circles[26].

The contrast between the two approaches extends beyond predictive performance. Classical rational models, such as those based on Nash equilibrium, provide a clear theoretical rationale for their predictions: they specify the incentives and equilibria

that should arise from rational choice. However, their explanatory transparency does not translate into empirical accuracy: Our results show that they systematically fail to reproduce the cooperation patterns observed in human data. LLM-based models, in contrast, may better approximate these empirical regularities, but they do so as black boxes whose internal reasoning remains inaccessible. Their predictions reveal what humans are likely to do, but not why. Nevertheless, unlike traditional models, LLMs can be directly interrogated through prompting and probing, offering a route to uncovering the latent heuristics that drive their choices. In this sense, LLMs complement rather than replace rational models: the former capture behavioral regularities, while the latter articulate mechanistic hypotheses. A deeper synthesis between interpretability and predictive validity may ultimately lead to a more complete understanding of social decision-making.

The distinct cooperation profiles across models provide additional insights. Qwen's behavior closely aligns with the Nash equilibrium (Table 2), which may reflect either its exposure to more game-theoretic literature during training or superior logical reasoning about optimal strategies. Interestingly, we deliberately avoided using game theory terminology such as "cooperate" and "defect" in our prompts to prevent triggering memorized responses, yet Qwen still converges toward rational play. While this makes Qwen an excellent test of strategic reasoning, it may render it less useful for predicting human behavior in novel games, since humans systematically deviate from Nash equilibrium predictions (Panel A of Figure 2). Mistral exhibits an intermediate pattern, showing sensitivity to the relative payoff structure but not fully optimizing in the game-theoretic sense.

To address the concern that LLMs might merely reproduce memorized information about the original experiments, we note two points. First, the novel games in our extended parameter space could not have been memorized, as they were never conducted with humans. Second, the substantial differences in cooperation profiles between models (even though all were trained on similar internet-scale corpora) suggest that their behaviors more likely reflect architectural and training differences than simple retrieval of stored experimental results.

A significant methodological contribution of our work is demonstrating that persona-based prompting is not necessary to replicate aggregate behavioral patterns. Rather than attempting to simulate individual participants, we show that focusing on population-level patterns provides a more tractable and robust approach. This finding simplifies the simulation process while maintaining predictive validity for the collective behavioral phenomena of interest.

The progressive improvement in results through our answer extraction techniques deserves elaboration. Simple extraction often produced nearly random responses because, when models are asked to provide a one-word answer, they default to surface-level pattern matching. By requesting detailed explanations (Double Extraction), we leveraged the phenomenon whereby models perform better when allowed to "think aloud" by articulating their reasoning to improve their logical consistency[27]. The Multi-step Prompt further decomposed complex reasoning into manageable steps, reducing complexity and thereby minimizing errors. Finally, the Logical Verifier acts as a quality control mechanism, filtering out responses containing mathematical errors or logical contradictions. This layered approach functions as a form of "attention check" for LLMs, ensuring that recorded decisions reflect genuine reasoning rather than random or inconsistent outputs. Together, these techniques help transform LLMs from unreliable proxies of human respondents into robust digital twins.

An intriguing anomaly appears in Llama's behavior at the extreme parameter values near (0,0), where cooperation appears minimal or random-like. This may reflect numerical instability when payoffs approach zero, making the strategic structure of the game less salient to the model. Further investigation of edge cases in the parameter space could reveal the boundaries of LLMs' strategic reasoning capabilities.

Importantly, the validity of our approach can only be definitively assessed through future experiments with human participants. While the precise mechanisms that enable LLMs to capture patterns of human experience and behavior may remain opaque and debated, we can evaluate the models' predictive power in a straightforward way: by carrying out subsequent experiments with humans, ideally following strictly specified experimental setups. To enable unambiguous evaluation, we have publicly pre-registered[1] the complete experimental design that needs to be conducted with human participants to validate our extended simulation results. This pre-registration specifies all details of the experimental protocol, preventing post-hoc adjustments and ensuring rigorous hypothesis testing.

More broadly, our work exemplifies the potential of AI for scientific discovery in the social and behavioral sciences. By using LLMs as computational models of human behavior, we can rapidly explore vast experimental spaces that would be prohibitively expensive or too time-consuming to investigate with human participants alone. This approach follows a virtuous cycle: validate models on existing human data, use models to generate novel hypotheses, then conduct targeted human experiments to test those hypotheses. The transparent, pre-registered nature of this process ensures that AI-generated predictions face genuine empirical tests, maintaining the rigor of scientific inquiry while dramatically expanding its scope and efficiency. As LLMs continue to improve, this methodology can be expected to become a standard tool for behavioral and social scientists, enabling systematic in-silico exploration of experimental landscapes before committing resources to the actual experiments involving humans.

# 3 Methods

## 3.1 Experimental Design and Models

As in the original experiment, we kept $R$ and $P$ constant at 10 and 5, respectively. To extend the parameter space, we varied $S$ and $T$ more broadly: $S$ ranges from 0 to 20, with 21 possible values, and $T$ ranges from 0 to 20, with 21 possible values. In total, we have simulated 441 games, 320 more than in the original experiments. Each game was played 20 times to generate sufficient data for identifying behavioral patterns. Our objective was to construct a cooperation matrix in which each element represents a specific game defined by its $(S, T)$ parameter combination, with the element's value indicating the average cooperation rate for that game.

## 3.2 Prompt engineering and model parameters

Model prompts consist of three different components (for the full text of prompts see Subsection A.2 in the Appendix):

- **System prompt.** Specifies the role assigned to the model. We use the default "You are a helpful assistant".

- **User prompt.** Contains the request made to the model. This prompt includes the game instructions and requests a choice.

- **Assistant prompt.** Contains the model's response. This component can be used to simulate a conversation history of previous exchanges.

All simulation code[2] was written in Python, and we used the vLLM library for efficient inference with the LLMs. Since each model requires a different prompting structure, we developed functions to adapt the general prompt structure to the format required by each LLM.

We adjusted the *temperature* and *max_tokens* parameters depending on the nature of the desired response. High *temperature* values promote diverse outputs but increase randomness, while low temperature produces more deterministic responses. The *max_tokens* parameter sets a maximum token limit for the output. The exact values for these parameters are specified in the relevant Subsections. We conducted our experiments on MareNostrum 5 at the Barcelona Supercomputing Center (BSC). For the main results in Panel B of Figure 2, we used 2 jobs, each running 9 hours and 20 minutes on one node using 20 Intel Sapphire Rapids 8460Y+ CPU cores, 1 H100 GPU with 64GB HBM2 memory and up to 134GB RAM.

## 3.3 Adapting the original game instructions to LLMs

We initially attempted to directly use the original instructions from the human experiment with minimal modification. For each simulated game round, we recorded the model's choice as 1 if it chose to cooperate (choice A or B, depending on random assignment) and 0 if it chose to defect. After completing all 20 iterations of each game, we calculated the average cooperation rate per game by averaging these binary outcomes.

After running initial test simulations with the original game instructions, we observed that the resulting prompts were overly verbose and that models prioritized generating human-like conversational responses rather than engaging in strategic reasoning about payoff structures. We therefore conducted an iterative refinement process, identifying the following key improvements:

- **Deletion of non-essential information.** We removed all contextual information about the experimental motivation, participating institutions, and background details. We retained only the core rules of the game and the payoff outcomes structure.

- **Emphasizing one-shot nature and simultaneity.** We explicitly stated that each game consisted of a single round with simultaneous decision-making to prevent an observed tendency for models to adopt cooperative strategies aimed at influencing future rounds or the opponent's current choice.

- **Avoiding explicit game-theoretic terminology.** We observed that using game theory vocabulary such as "opponent", "cooperate" or "strategies" led models to activate memorized game-theoretic scenarios rather than reasoning about the specific payoff structure. We therefore adopted neutral vocabulary such as "other player" instead of "opponent" and replaced C and D with A and B as the two choice options.

- **Clarifying no competitive objective.** We observed that models often treated the games as zero-sum competitions, choosing to defect to "beat" the other player even in Harmony Games where this was disadvantageous for both parties. We therefore explicitly stated that each player's prize depends solely on their own points earned, not on relative performance or point differences.

---

[2]https://github.com/acerapal/Replicating-Human-Game-Theory-Experiments-with-LLMs-Phenotypic-Profiles-and-Bias-Modulation

- **Clarifying the prize structure.** We replaced the original lottery ticket rewards with a clearer monetary conversion: 10 euros per point. We also added a concrete example of the calculation (e.g., "if you earn 8 points, you receive 80 euros") to ensure models understood the direct relationship between points and payoffs and to reduce spurious calculations.

These changes were developed through manual prompt engineering informed by established best-practice guides[28, 29]. Our iterative refinement process consisted of: (1) running simulations with the current prompt, (2) analyzing model responses across multiple games to identify systematic errors or misconceptions, (3) modifying the prompt to address these issues, and (4) validating that the problematic patterns were resolved in subsequent simulations. This cycle was repeated until the models provided plausible answers that were consistent with the game structure and payoff incentives.

## 3.4 Answer extraction approaches
We developed a four-stage extraction methodology with progressively increasing complexity to reliably extract and validate strategic choices from model outputs.

### 3.4.1 Generation of long answers
Initially, we prompted models for brief responses that could be easily parsed using regular expressions. However, these short answers exhibited high variability and often appeared arbitrary. In contrast, when we allowed models to generate longer responses, they produced more logically coherent reasoning. This suggests that having the space to articulate explanations helps the models perform, a phenomenon observed in prior work showing that encouraging models to "think step-by-step" enhances task accuracy[30].

Based on these observations, we adopted a Double Extraction approach: (1) elicit a detailed response from the tested model, then (2) use a second LLM to extract the final choice from this explanation. Long answers were generated by providing the instruction prompt (Subsection A.2 in the Appendix) with the game-specific payoff values. To prevent labeling bias, we randomly assigned the labels A and B to the strategies C (cooperate) and D (defect) across games. For generating long answers, we used a *temperature* of 0.8 (to encourage diverse reasoning) and set *max_tokens* to 1000.

### 3.4.2 Extraction of short answers
The second step of the Double Extraction process extracts the final choice (in a few words only) from the detailed response. We provided the long answer generated by the tested LLM as input to a second LLM (Qwen), which we prompted to identify only the player's choice: A or B. A simple regular expression function then parsed this output to extract the selected letter. If the extraction successfully identified A or B, we recorded the choice as valid; otherwise, we flagged it as invalid and required the tested model to replay that game.

| Qwen | Llama | Mistral |
|---|---|---|
| 0.97 | 0.96 | 0.83 |

**Table 3. Extraction accuracy for identifying player choices from long-form responses.** Accuracy rates of each model when tasked with extracting the final choice (A or B) from detailed reasoning outputs, validated against manually annotated ground truth data.

To select the most reliable model for short answer extraction, we conducted a manual validation study. We randomly sampled 100 long answers (distributed equally across the three tested LLMs) and manually annotated each with the correct choice: A, B, or neither. We then used each of the three models (Llama, Mistral, and Qwen) as extractors and calculated their accuracy rates. Table 3 shows that Llama and Qwen performed comparably as extractors (0.96 and 0.97 accuracy, respectively), while Mistral lagged substantially behind (0.83). We selected Qwen for all subsequent extractions due to its marginally higher accuracy and to maintain consistency with its role in the logical verification step (described below). For the extraction task, we used a lower *temperature* of 0.3 to prioritize deterministic, concise outputs over creative elaboration, and set *max_tokens* to 50.

### 3.4.3 Multi-step prompt
Although the double extraction approach yielded improved results compared to simple extraction, manual inspection revealed that the models still produced frequent logical inconsistencies and factually incorrect statements. Previous research[27, 30] has shown how results can be improved by prompting models to articulate step-by-step reasoning. Breaking down complex reasoning tasks into manageable steps has also been shown to be helpful[28, 29]. We therefore designed prompts that explicitly guided models to decompose their reasoning into discrete steps.

These steps needed to be specific enough to scaffold the reasoning process, yet generic enough to avoid biasing models toward particular strategies. We iteratively tested multiple multi-step prompt variants, each of which instructed models to: (1)

group the four possible outcomes by their own choice (A or B), and (2) compare the payoffs within each group. To identify the optimal prompt, we evaluated performance in the $S \geq T$ region (Harmony Games), where cooperation is the only sensible choice. We selected the prompt version that yielded the highest average cooperation rates for both Llama and Mistral in this region. This multi-step instruction was inserted into the final prompt immediately following the description of game outcomes.

### 3.4.4 Logical verifier

Despite the multi-step prompt improvements, manual inspection continued to reveal logical and mathematical errors in model responses. We therefore implemented a logical verification layer as a final quality control mechanism.

We first conducted systematic error analysis, cataloging the most common mistakes produced by Llama and Mistral. Based on these patterns, we designed a verification prompt that instructed Qwen (maintaining consistency with the extraction step) to classify each long-form response as "good" or "bad" based on its logical validity. The prompt consisted of two components: First, we explicitly defined criteria for valid and invalid responses in bullet-point format, emphasizing the specific logical errors we had identified (e.g. incorrect arithmetic comparisons, misunderstanding of outcome probabilities or inconsistencies between reasoning and final choice). Second, we incorporated diverse examples of both acceptable and flawed responses, employing the few-shot learning technique[31], which has been shown to enhance logical reasoning even in mid-sized models like those used in our experiments. We tested multiple versions of this combined prompt with varying example sets and evaluated their performance through simulation. Following the same validation approach used for the multi-step prompt, we selected the version that maximized cooperation rates in the $S \geq T$ region (Harmony Games).

## 3.5 Conditions for answer validity

Each game required a valid response to be included in the final dataset. Our validation pipeline proceeded as follows: (1) the tested model generated a long-form response, (2) the logical verifier (Qwen) classified it as valid or invalid, (3) if valid, the extractor (also Qwen) identified the final choice (A or B). Games producing invalid responses at either verification or extraction stages were replayed in subsequent rounds until a valid response was obtained.

However, preliminary testing revealed that certain game configurations consistently produced logically flawed responses, trapping the simulation in infinite loops. To address this issue, we implemented an adaptive relaxation mechanism: if the number of invalid games remained unchanged between two consecutive rounds, we temporarily disabled the logical verifier for those problematic games, allowing them to proceed directly to the extraction step. This pragmatic compromise ensured simulation completion while maintaining quality control for the majority of responses (for details see 'Subsection A.4 in Appendix).

## 3.6 Computation of Nash equilibrium

To compute Nash equilibrium cooperation rates, we simulated replicator dynamics from evolutionary game theory. When strategies compete, the fraction of strategy $i$ evolves according to:

$$\dot{x}_{i,t} = x_{i,t} \left( \pi_i(x_t) - \bar{\pi}(x_t) \right) \tag{1}$$

where $x_{i,t}$ represents the fraction of the strategy $i$ at time $t$, $\pi_i(x_t)$ is its average payoff and $\bar{\pi}(x_t)$ is the average payoff of all strategies at time t.

With only two strategies (C and D), we just need to track the cooperator fraction. Setting $\Delta t = 1$ in Equation 1, we obtain:

$$x_{t+1} = x_t(1-x_t)\left(\pi_C(x_t) - \pi_D(x_t)\right) + x_t \tag{2}$$

where $x_t$ is cooperator fraction at time $t$, $\pi_C(x_t)$ and $\pi_D(x_t)$ are cooperator and defector average payoffs, respectively. These payoffs are:

$$\pi_C = x_t R + (1-x_t)S \tag{3}$$
$$\pi_D = x_t T + (1-x_t)P \tag{4}$$

with $R = 10$ and $P = 5$ as in the original experiments[22], while $S$ and $T$ vary by game.

We iterated until $t = 10^3$ or until one of four outcomes occurred (tolerance $\varepsilon = 0.1$):

- **Total Cooperation:** $x_t \geq 1 - \varepsilon$ (Harmony Game).

- **Total Defection:** $x_t \leq \varepsilon$ (Prisoner's Dilemma).

- **Mixed Equilibrium:** $|\pi_C - \pi_D| \leq \varepsilon$, where payoffs equalize and cooperator fraction stabilizes between 0 and 1 (Snowdrift Game, Stag Hunt diagonal).

- **Periodicity:** Oscillation between two $x_t$ values; we recorded their average (some Snowdrift games).

We have verified these results by performing a fixed-point and stability analysis of the replicator dynamics corresponding to the $2 \times 2$ symmetric game defined by the payoff matrix:

$$\begin{pmatrix} R & S \\ T & P \end{pmatrix}$$

Here, the two strategies correspond to *Cooperate* (C) and *Defect* (D). Let $x \in [0, 1]$ denote the fraction of cooperators in the population. Using Equation 1, the replicator dynamics can be written as

$$\dot{x} = x(1-x)\left[\pi_C(x) - \pi_D(x)\right] \tag{5}$$

where

$$\pi_C(x) = Rx + S(1-x) \tag{6}$$
$$\pi_D(x) = Tx + P(1-x) \tag{7}$$

Equation 5 can be rewritten in the compact form

$$\dot{x} = x(1-x)g(x) \tag{8}$$

where

$$g(x) = \pi_C(x) - \pi_D(x) = (R - T - S + P)x + (S - P) \tag{9}$$

The function $g(x)$ measures the instantaneous advantage of cooperation over defection: if $g(x) > 0$, cooperation tends to increase; if $g(x) < 0$, it tends to decrease. Fixed points correspond to values of $x$ for which $\dot{x} = 0$, i.e.

$$x \in \{0, 1\} \quad \text{or} \quad g(x) = 0.$$

**(i) Boundary fixed points.** To test the stability of the homogeneous states $x = 0$ (all defect) and $x = 1$ (all cooperate), we can inspect the sign of $\dot{x}$ in their neighborhoods:

- Near $x = 0$ (almost all defect), $x$ is small and the factor $(1-x) \approx 1$, so

$$\dot{x} \approx x g(0), \quad \text{where} \quad g(0) = S - P.$$

If $g(0) > 0$ ($S > P$), then $\dot{x} > 0$ for small $x$: cooperation increases and $x = 0$ is *unstable*. Conversely, if $g(0) < 0$ ($S < P$), then $\dot{x} < 0$ near $x = 0$ and the population returns to full defection: $x = 0$ is *stable*.

- Near $x = 1$ (almost all cooperate), set $x = 1 - \varepsilon$ with $\varepsilon \ll 1$. Then

$$\dot{x} \approx \varepsilon g(1), \quad \text{where} \quad g(1) = R - T.$$

If $g(1) > 0$ ($R > T$), then $\dot{x} > 0$ for $x < 1$, so the system moves back toward $x = 1$: cooperation is *stable*. If $g(1) < 0$ ($R < T$), then $\dot{x} < 0$ and $x$ decreases away from 1: cooperation is *unstable*.

These two simple conditions,

$$x = 0 \text{ defection stable if } S < P, \qquad x = 1 \text{ cooperation stable if } R > T,$$

partition the parameter space into regions of dominance and coordination.

**(ii) Interior fixed point.** An interior equilibrium $x^*$ exists whenever $g(x^*) = 0$, that is,

$$x^* = \frac{P-S}{R-T-S+P}. \tag{10}$$

The interior fixed point exists if $(0 \leq x^* \leq 1)$ and $(R-T-S+P) \neq 0$. Let us rewrite its expression as

$$x^* = \frac{P-S}{D}, \qquad D := R-T-S+P.$$

Then

$$x^* > 0 \iff (P-S)D > 0 \quad \text{and} \quad x^* < 1 \iff \frac{R-T}{D} > 0 \iff (R-T)D > 0,$$

since $1-x^* = \frac{R-T}{D}$. Therefore,

$$0 < x^* < 1 \iff (P-S)D > 0 \text{ and } (R-T)D > 0,$$

i.e., $D$ must have the *same sign* as both $P-S$ and $R-T$. Therefore, we have two regions of existence:

$$\begin{cases} D > 0 \Rightarrow S < P \text{ and } T < R \\ D < 0 \Rightarrow S > P \text{ and } T > R \end{cases}$$

*Specialization to our parameters.* With $R = 10$ and $P = 5$,

$$0 < x^* < 1 \iff (S < 5 \text{ and } T < 10) \text{ or } (S > 5 \text{ and } T > 10),$$

corresponding exactly to the bottom-left and top-right regions of the $(T,S)$ plane.

To test its stability, we can linearize Equation 8 around $x^*$:

$$\dot{x} \approx x^*(1-x^*)g'(x^*)(x-x^*),$$

since $g(x^*) = 0$. The factor $x^*(1-x^*) > 0$, so the sign of $g'(x^*)$ determines the stability: Stable if $g'(x^*) < 0$, unstable if $g'(x^*) > 0$. From Equation 9,

$$g'(x) = R-T-S+P = D,$$

which is constant. Therefore $x^*$ is stable if $D = R-T-S+P < 0$, which corresponds to the top right quadrant of our region of study ($S > 5$ and $T > 10$).

Intuitively, this condition means that when the gain from mutual cooperation $(R-P)$ is smaller than the temptation and sucker effects combined $(T-S)$, the dynamics settle into a stable coexistence between cooperators and defectors (anti-coordination). Otherwise, the dynamics are bistable, typical of coordination games.

**(iii) Summary of Stability Criteria.**

| Condition | Stable point(s) | Game type |
|---|---|---|
| $S < P, T > R$ (bottom right quadrant) | $x = 0$ | Defection dominates (Prisoner's Dilemma) |
| $S > P, T < R$ (top left quadrant) | $x = 1$ | Cooperation dominates (Harmony) |
| $S < P, R > T, R-T-S+P > 0$ (bottom left quadrant) | $x = 0, 1$ (bistable) | Coordination (Stag Hunt) |
| $S > P, R < T, R-T-S+P < 0$ (top right quadrant) | $x^*$ (mixed) | Anti-coordination (Snowdrift) |

In the coordination (bistable) region, corresponding to $T < 10$ and $S < 5$, both $x = 0$ (full defection) and $x = 1$ (full cooperation) are locally stable, while the interior fixed point $x^*$ is unstable. Since in our simulations we set $x(0) = 0.5$, the long-term outcome depends on the position of this initial condition relative to $x^*$:

$$\begin{cases} x(0) > x^* \implies \text{the population converges to full cooperation } (x = 1), \\ x(0) < x^* \implies \text{the population converges to full defection } (x = 0). \end{cases}$$

Using Equation (10), this boundary between the basins of attraction is given by

$$x^* = \frac{5 - S}{15 - T - S},$$

so that for $x(0) = 0.5$ cooperation prevails whenever $T < 5 + S$, and defection otherwise. This condition precisely separates the two attraction basins observed in the numerical simulations.

# References

1. Argyle, L. P. *et al.* Out of One, Many: Using Language Models to Simulate Human Samples. *Polit. Analysis* **31**, 1–15, DOI: 10.1017/pan.2023.2 (2023).

2. Rahwan, I. *et al.* Machine behaviour. *Nature* **568**, 477–486, DOI: 10.1038/s41586-019-1138-y (2019).

3. Meng, X. *et al.* The application of large language models in medicine: A scoping review. *iScience* **27**, 109713–109713, DOI: 10.1016/j.isci.2024.109713 (2024).

4. Thirunavukarasu, A. J. *et al.* Large language models in medicine. *Nat. Med.* **29**, 1930–1940, DOI: 10.1038/s41591-023-02448-8 (2023).

5. Lai, J. *et al.* Large language models in law: A survey (2023). 2312.03718.

6. Homoki, P. & Ződi, Z. Large language models and their possible uses in law. *Hungarian J. Leg. Stud.* **64**, 435–455, DOI: 10.1556/2052.2023.00475 (2024).

7. Wang, S. *et al.* Large language models for education: A survey and outlook (2024). 2403.02334.

8. Grassucci, E. *et al.* Beyond answers: How llms can pursue strategic thinking in education (2025). 2504.04815.

9. Pasupuleti, M. K. *AI's role in global stability, diplomacy, and peacebuilding*, 1–11 (National Education Services, 2025).

10. Binz, M. *et al.* A foundation model to predict and capture human cognition. *Nature* **644**, 1002–1009, DOI: 10.1038/s41586-025-09215-4 (2025).

11. Xie, H. & Zhu, J.-Q. Centaur May Have Learned a Shortcut that Explains Away Psychological Tasks, DOI: 10.31234/osf.io/u7z4t_v1 (2025).

12. Sturm, B. L. A Simple Method to Determine if a Music Information Retrieval System is a "Horse". *IEEE Transactions on Multimed.* **16**, 1636–1644, DOI: 10.1109/TMM.2014.2330697 (2014).

13. Horton, J. J. Large language models as simulated economic agents: What can we learn from homo silicus? (2023). 2301.07543.

14. Aher, G., Arriaga, R. I. & Kalai, A. T. Using large language models to simulate multiple humans and replicate human subject studies. In *Proceedings of the International Conference on Machine Learning (ICML)*, DOI: 10.48550/arXiv.2208.10264 (2023).

15. Mei, Q., Xie, Y., Yuan, W. & Jackson, M. O. A turing test of whether ai chatbots are behaviorally similar to humans. *Proc. Natl. Acad. Sci. USA* **121**, DOI: 10.1073/pnas.2313925121 (2024).

16. Hagendorff, T. *et al.* Machine Psychology, DOI: 10.48550/arXiv.2303.13988 (2024). 2303.13988.

17. Pellert, M., Lechner, C. M., Wagner, C., Rammstedt, B. & Strohmaier, M. AI Psychometrics: Assessing the Psychological Profiles of Large Language Models Through Psychometric Inventories. *Perspectives on Psychol. Sci.* DOI: 10.1177/17456916231214460 (2024).

18. Grattafiori, A. *et al.* The Llama 3 Herd of Models, DOI: 10.48550/arXiv.2407.21783 (2024). 2407.21783.

19. Jiang, A. Q. *et al.* Mistral 7B, DOI: 10.48550/ARXIV.2310.06825 (2023).

20. Yang, A. *et al.* Qwen2 Technical Report, DOI: 10.48550/arXiv.2407.10671 (2024). 2407.10671.

21. Qwen *et al.* Qwen2.5 Technical Report, DOI: 10.48550/arXiv.2412.15115 (2025). 2412.15115.

22. Poncela-Casasnovas, J. *et al.* Humans display a reduced set of consistent behavioral phenotypes in dyadic games. *Sci. Adv.* **2**, e1600451, DOI: 10.1126/sciadv.1600451 (2016).

23. Simon, H. A. A Behavioral Model of Rational Choice. *The Q. J. Econ.* **69**, 99–118, DOI: 10.2307/1884852 (1955). 1884852.

24. Feldman, J. The simplicity principle in perception and cognition. *WIREs Cogn. Sci.* **7**, 330–340, DOI: 10.1002/wcs.1406 (2016).

25. Martin-Gutierrez, S., van Dissel, M. N. C. & Karimi, F. The hidden architecture of connections: How do multidimensional identities shape our social networks?, DOI: 10.48550/arXiv.2406.17043 (2024). 2406.17043.

26. Tamarit, I., Cuesta, J. A., Dunbar, R. I. M. & Sánchez, A. Cognitive resource allocation determines the organization of personal networks. *Proc. Natl. Acad. Sci.* **115**, 8316–8321, DOI: 10.1073/pnas.1719233115 (2018).

27. Wei, J. *et al.* Chain-of-thought prompting elicits reasoning in large language models (2022). 2201.11903.

28. Fagbohun, O., Harrison, R. M. & Dereventsov, A. An empirical categorization of prompting techniques for large language models: A practitioner's guide. *J. Artif. Intell. Mach. Learn. Data Sci.* **1**, 1–11, DOI: 10.51219/jaimld/oluwole-fagbohun/15 (2023).

29. Amatriain, X. Prompt design and engineering: Introduction and advanced methods (2024). 2401.14423.

30. Kojima, T. *et al.* Large language models are zero-shot reasoners. In *Advances in Neural Information Processing Systems (NeurIPS)*, DOI: 10.48550/arXiv.2205.11916 (2022).

31. Brown, T. B. *et al.* Language models are few-shot learners. In *Advances in Neural Information Processing Systems*, vol. 33, 1877–1901, DOI: 10.48550/arXiv.2005.14165 (2020).

32. Kahneman, D., Knetsch, J. L. & Thaler, R. H. Fairness as a constraint on profit seeking: Entitlements in the market. *Amer. Econ. Rev.* **76**, 728–741 (1986).

33. Güth, W., Schmittberger, R. & Schwarze, B. An experimental analysis of ultimatum bargaining. *J. Econ. Behav. Organ.* **3**, 367–388, DOI: 10.1016/0167-2681(82)90011-7 (1982).

34. Berg, J., Dickhaut, J. & McCabe, K. Trust, reciprocity, and social history. *Games Econ. Behav.* **10**, 122–142, DOI: 10.1006/game.1995.1027 (1995).

35. Rapoport, A. & Chammah, A. M. *Prisoner's Dilemma* (University of Michigan Press, 1965).

36. Hong, L. & Page, S. E. Groups of diverse problem solvers can outperform groups of high-ability problem solvers. *Proc. Natl. Acad. Sci. USA* **101**, 16385–16389, DOI: 10.1073/pnas.0403723101 (2004).

37. Milgram, S. Behavioral study of obedience. *J. Abnorm. Soc. Psychol.* **67**, 371–378, DOI: 10.1037/h0040525 (1963).

## Author contributions statement

ACP: Conceptualization, Methodology, Formal Analysis, Data Curation, Visualization, Writing - Original Draft, Writing - Review and Editing. SMG: Original Conceptualization, Formal Analysis, Writing - Review and Editing. ASC: Conceptualization, Writing - Review and Editing. MP: Original Conceptualization, Methodology, Visualization, Writing - Original Draft, Writing - Review and Editing. All of the authors approved the final manuscript for submission.

## Additional information

Code and data to fully replicate our analyses is available at https://github.com/acerapal/Replicating-Human-Game-Theory-Experiments-with-LLMs-Phenotypic-Profiles-and-Bias-Modulation.

# A Appendix

## A.1 Details on games and behavioral experiments

- **Dictator game.** One player is given an amount of money and can decide how much to keep and how much to give to the other player. Humans tend to keep most of the money[32]. This game is used to measure altruism.

- **Ultimatum game.** This game is similar to the Dictator game, but the recipient can reject the offer, resulting in neither player receiving anything. In humans,[33] the allocations are more balanced than in the Dictator game due to fear of punishment. This game is used to measure fairness, altruism, and responses to unfairness.

- **Trust game.** One player, the investor, receives a quantity of money and decides how much to keep and how much to give to the other player, the banker. The banker multiplies the received money by a factor and decides how much to return to the investor. Humans[34] tend to show reciprocity. This game measures altruism and fairness.

- **Prisoner's dilemma.** This game consists of two players who can either cooperate or defect. If one player defects while the other cooperates, they receive the best outcome, and the other, the worst. If both players cooperate, they receive a higher joint payoff. Humans[35] show a tendency toward cooperation and are influenced by past outcomes. This game measures collective well-being versus personal interest.

- **Wisdom of crowds.** An observed phenomenon whereby averaging the responses to a question yields more accurate results than individual responses[36]. This phenomenon measures collective intelligence.

- **Milgram obedience experiment.** This experiment involves an authority figure, a participant who must obey, and a subject who experiences pain (simulated electric shocks) whenever the authority figure dictates, administered by the obedient participant. It measures obedience to authority. In humans[37], obedience remains high until the subject receiving pain begins to complain.

## A.2 Final Prompts

We use three different types of prompts:

- **Instruction Prompt:** General instructions of the game and explanation of the different payoffs.

- **Multi-steps Prompt:** Decomposition into distinct logical steps to facilitate the models' reasoning.

- **Logical verifier Prompt:** Definition of good and bad long answers, examples to facilitate the classification and a request for output (i.e. the classification).

**Instruction Prompt with exemplary payoff values** .

```
This one-shot game, is made of only one round with an anonymous player. You will play only
    once.

To play you must choose one of two options: A and B, the other player will also choose
    between A and B. Both players are subjected to the same rules.

You won't know the other player decision until the end of the round so you won't be able to
    change your choice after knowing the other player's decision.
The other player won't know your decision until the end of the round so he won't be able to
    change his choice after knowing yours.

You will be playing simultaneously with the other player.

Both players will receive a monetary prize. The monetary prize for the player will be the
    amount of points the player has earned multiplied by 10.

For example, if player 1 earns 9 points and player 2 earns 5 points.
Player 1 will receive 9*10 = 90 euros and player 2 will receive 10*5= 50 euros.

The monetary prize of each player only depends on his number of points.
The difference in points of both players has no effect on the prize.
```

```
Both players will receive the prize, regardless of who earned more or less money.
The outcome is solely determined by the number of points earned by each player.


If you choose A and the other player chooses A. You earn 10 points, the other player earns
    10 points.
If you choose A and the other player chooses B. You earn 6 points, the other player earns 11
     points.
If you choose B and the other player chooses A. You earn 11 points, the other player earns 6
     points.
If you choose B and the other player chooses B. You earn 5 points, the other player earns 5
    points.
```

**Multi-step Prompt** .

```
1. There are four possible outcomes in this game,
depending on your choice (A or B) and the other player's choice (A or B).
2. Group the outcomes based on your decision:
   – If you choose A: (A,A) and (A,B)
   – If you choose B: (B,A) and (B,B)
3. Analyze and compare the outcomes in each group step by step.
4. Make sure your analysis is consistent with the rules of the game
and the possible outcomes.
5. Based on your analysis, choose either A or B.
```

**Logical verifier prompt** .

```
    You have to read the answer of Player 1 and output good or bad
according to the following criteria:

Criteria for good:
 – ALL arithmetic/mathematic comparisons and claims are CORRECT.
 – ALL descriptions of the possible outcomes are CORRECT.
 – ONLY TWO possible outputs per choice, for A: (A,A) and (A,B), for B: (B,A) and (B,B).
 – The final choice is CONSISTENT with the previous reasoning.
 – ALL statements make sense.
 – GOOD UNDERSTANDING of the rules.

Criteria for bad:
 – ONE OR MORE arithmetic/mathematic comparisons or claims are WRONG.
 – ONE OR MORE descriptions of the possible outcomes are WRONG.
 – Does not understand that there are ONLY TWO possible outcomes per choice.
 – Final choice is NOT CONSISTENT with the previous reasoning.
 – ONE OR MORE statements do not make sense.
 – BAD UNDERSTANDING of the rules.

Rules of the game:
{instructions_script_short}

Example 1, bad reasoning:

Rules of Example 1:

If you choose A and the other player chooses A. You earn 10 points, the other player earns
    10 points.
If you choose A and the other player chooses B. You earn 7 points, the other player earns 6
    points.
If you choose B and the other player chooses A.
You earn 6 points, the other player earns 7 points.
If you choose B and the other player chooses B. You earn 5 points, the other player earns 5
    points.
```

```
Answer of Example 1:

Based on the rules and possible outcomes, here's the analysis:

- If you choose A, you can earn either 10 points (if the other player chooses A) or 7 points
     (if the other player chooses B).
- If you choose B, you can earn either 5 points (if the other player chooses B) or 6 points
    (if the other player chooses A).

However, the rules state that the difference in points between you and the other player has
    no effect on the prize.
Therefore, it comes down to the following:

- If you choose A, the highest possible prize is 10 points, but the lowest possible prize is
     7 points.
- If you choose B, the highest possible prize is 5 points, but the lowest possible prize is
    6 points.

Since the goal is to maximize the possible prize, it would seem that choosing A gives you
    the potential for a higher prize (100 euros) compared to choosing B (60 euros as the
    highest potential prize).
However, it's important to keep in mind that the other player is making the same analysis as
     well, so it might be wiser to choose B to reduce the chances of both players choosing A
and thus reducing the potential prize for both.

As this is a one-shot game, it's difficult to say for certain which choice will lead to a
    higher average prize for both players. It could be a game of bluff and strategy rather
    than a purely mathematical decision.

In this analysis, I would lean towards choosing B, to potentially lower the chances of both
    players receiving a lower prize, but it's ultimately up to your personal decision.

Output of Example 1: bad

Example 2, good reasoning:

Rules of Example 2:

If you choose A and the other player chooses A. You earn 10 points, the other player earns
    10 points.
If you choose A and the other player chooses B. You earn 10 points, the other player earns 6
     points.
If you choose B and the other player chooses A. You earn 6 points, the other player earns 10
     points.
If you choose B and the other player chooses B. You earn 5 points, the other player earns 5
    points.

Answer of Example 2:

To make a decision, let's analyze the outcomes based on the rules and possible outcomes:

1. If you choose A:
    - (A,A): Both players earn 10 points each.
    - (A,B): You earn 10 points, and the other player earns 6 points.

2. If you choose B:
    - (B,A): You earn 6 points, and the other player earns 10 points.
    - (B,B): Both players earn 5 points each.
```

3. Since the monetary prize for each player is the number of points multiplied by 10, it is
   clear that earning more points results in a higher prize.

4. Comparing the outcomes in each group:
   - If you choose A: (A,A) gives you the highest reward, and (A,B) gives you the same
     reward but a lower-one for the other player.
   - If you choose B: (B,B) gives you the lowest reward, and (B,A) gives a slightly higher
     reward, but still lower than all possible rewards when choosing A.

5. Based on this information, choosing A would result in a higher potential reward, making
   it the better choice.
Even though one might worry about the other player choice, it doesn't make sense since the
   only possible outcome of choosing A is 10 points, which is the maximum.
The other player choice isn't going to affect your outcome at all, there is no risk.

Output of Example 2: good

Example 3, bad reasoning:

Rules of Example 3:

If you choose A and the other player chooses A. You earn 10 points, the other player earns
   10 points.
If you choose A and the other player chooses B. You earn 9 points, the other player earns 15
    points.
If you choose B and the other player chooses A. You earn 15 points, the other player earns 9
    points.
If you choose B and the other player chooses B. You earn 5 points, the other player earns 5
   points.

Answer of Example 3:

Let's analyze the outcomes:
- If you choose A, the probability of the other player also choosing A is 50%.
In this case, you would earn 10 points, and the other player would also earn 10 points (
   outcome A,A).
The other possibility is that the other player chooses B, in which case you would earn 9
   points, and they would earn 15 points (outcome A,B).
- If you choose B, the probability of the other player also choosing B is 25%.
In this case, you would earn 5 points, and the other player would also earn 5 points (
   outcome B,B).
The other possibility is that the other player chooses A, in which case you would earn 15
   points, and they would earn 9 points (outcome B,A).

Now, let's compare the outcomes:
- If you choose A, you have two possible outcomes: one with a higher reward (10 points) and
   one with a lower reward (9 points).
- If you choose B, you have two possible outcomes: one with a lower reward (5 points) and
   one with a higher reward (15 points).

Since you don't know the other player's choice, it's impossible to determine the exact
   outcome.
However, it's interesting to note that if the other player chooses A, choosing A would
   result in a lower reward (9 points) compared to choosing B (15 points).
On the other hand, if the other player chooses B, choosing A would result in a higher reward
    (10 points) compared to choosing B (5 points).
Given this analysis, it seems that choosing B gives you the potential for a higher reward,
   but it's also possible that the other player will choose A in which case you would have
   a lower reward.
It's truly a gamble. In this game, there is no optimal strategy as the game relies on the

other player's choice, which remains unknown.
Therefore, it's hard to make a definitive choice between A and B as there is no clear
    advantage in either option.
Based on the analysis and the lack of a perfect answer, I'll leave the choice up to you.
You can choose either A or B as there's no guaranteed outcome. Good luck!

Output of Example 3: bad

Example 4, good reasoning:

Rules of Example 4:

If you choose B and the other player chooses B. You earn 10 points, the other player earns
    10 points.
If you choose B and the other player chooses A. You earn 8 points, the other player earns 5
    points.
If you choose A and the other player chooses B. You earn 5 points, the other player earns 8
    points.
If you choose A and the other player chooses A. You earn 5 points, the other player earns 5
    points.

Answer of Example 4:

Based on the given information, I will analyze the outcomes for each choice:

If you choose A:
- If the other player chooses A, you both will receive 5 points each.
- If the other player chooses B, you will receive 5 points, while the other player receives
    8 points.

If you choose B:
- If the other player chooses A, you will receive 8 points, while the other player receives
    5 points.
- If the other player chooses B, you both will receive 10 points each.

If you choose B, the maximum amount of points you can get is 10 points and the minimum, 8
    points.
If you choose A the maximum amount of points you can get is 5 points and the minimum, 5
    points.
The maximum points is greater when you choose B, leading to a higher prize.
The minimum points is greater when you choose B, leading to a lower risk of getting little
    points.
Then, choice B is optimal since it is the option with the highest reward and the lower risk.

Output of Example 4: good

Example 5, bad reasoning:

Rules of Example 5:

If you choose B and the other player chooses B. You earn 10 points, the other player earns
    10 points.
If you choose B and the other player chooses A. You earn 10 points, the other player earns 7
     points.
If you choose A and the other player chooses B. You earn 6 points, the other player earns 10
     points.
If you choose A and the other player chooses A. You earn 5 points, the other player earns 5
    points.

Answer of Example 5:

```
Given the rules, choosing A would bring a maximum of 5 points, and choosing B would bring a
    maximum of 10 points if the other player also chooses B.
Since we don't know the other player's choice, it's impossible to guarantee a higher score
    by choosing A or B.
However, if we're forced to make a choice, choosing B seems like the better option since it
    has a higher potential reward if the other player also chooses B.
But, ultimately, this game is a gamble, and both choices have an element of risk.

Outcome of Example 5: bad

Example 6, good reasoning:

Rules of Example 6:

If you choose A and the other player chooses A. You earn 10 points, the other player earns
    10 points.
If you choose A and the other player chooses B. You earn 1 points, the other player earns 15
    points.
If you choose B and the other player chooses A. You earn 15 points, the other player earns 1
    points.
If you choose B and the other player chooses B. You earn 5 points, the other player earns 5
    points.

Answer of Example 6:

Choosing A can lead to the worst result (earning only 1 point) and choosing B can lead to
    the best result (earning 15 points).
Then I choose B to earn at least 5 points and also have the potential to the highest reward.

Output of Example 6: good

Rules of game to analyze:
{points}

Answer of Player 1, to analyze:
{answer1}

Important:
- Output ONLY one word: good or bad
- Do not add punctuation, extra spaces, or explanations.
```

## A.3  Additional Models (Qwen, Mistral)

While the main analysis focuses on results obtained with Llama, we also conducted the same simulations using the other models described earlier: Qwen and Mistral. As shown in Figure A1, progressively more sophisticated extraction methods reveal clearer cooperation patterns in both models, mirroring the improvements observed with Llama. Qwen exhibits a non-random cooperation matrix even with Simple Extraction, and Double Extraction is sufficient to obtain a clear pattern. In contrast, Mistral does not display a distinguishable cooperation matrix until the Multi-steps extraction stage, demonstrating greater sensitivity to the extraction methodology.

## A.4  Problematic Games

When the logical verifier cannot accept any answer in an iteration, it deactivates for that game's remaining responses, which proceed directly to answer extraction. Cooperation rates from these less-filtered games may reflect logical inconsistencies or mathematical errors rather than genuine strategic reasoning. Figure A2 shows the fraction of responses for each game that bypassed the logical verifier. The pattern reveals that games with higher values of $S$ and $T$ are more problematic, producing fewer logically correct answers and thus yielding cooperation estimates with potentially higher error. However, the overall impact remains limited: bypass rates do not exceed 0.25 for any game, indicating that the majority of responses successfully passed logical verification even in the most problematic parameter regions.
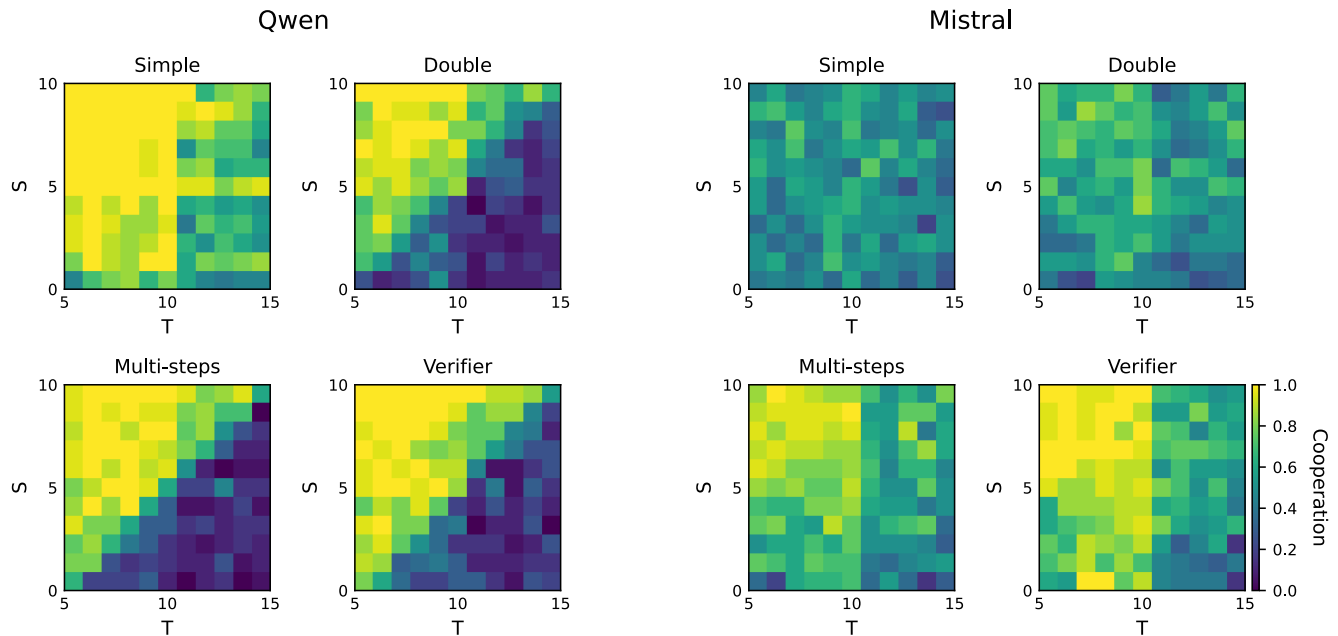
**Figure A1.** Progressive improvements in Mistral's and Qwen's cooperation matrices with increasingly sophisticated extraction methods. Each panel shows results from Simple Extraction, Double Extraction, Multi-step Extraction, and Logical Verifier Extraction (from top to bottom for each model). Color scales indicate average cooperation ranging from 0 (yellow: no cooperation) to 1 (purple: full cooperation). Note that Qwen exhibits clear patterns earlier in the extraction process, while Mistral requires more complex methods to reveal structured cooperation behavior.
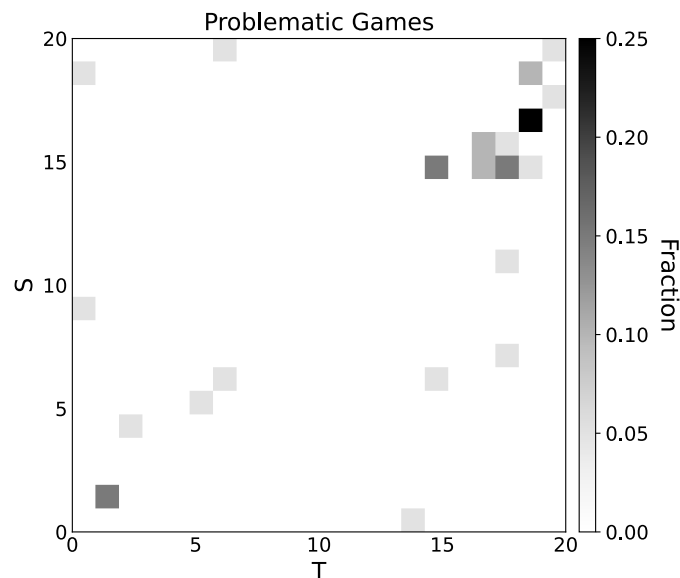


**Figure A2. Fraction of responses bypassing the logical verifier.** When at least one response for a game passes logical verification in an iteration, the verifier deactivates for that game's remaining responses, which proceed directly to answer extraction. This figure shows the fraction of responses that bypassed verification for each game, with values ranging from 0 (white: all responses accepted by logical verifier) to 0.25 (black: 1 in 4 responses bypassed logical verifier).