ISSN 2392-1641 e-ISSN 2450-0097

Economics and Business Review

Volume 11 (2) 2025

CONTENTS

Editorial introduction Joanna Lizińska, Katarzyna Schmidt-Jessa, Konrad Sobański

ARTICLES

How initial price history influences expectation formation in multi-asset experimental markets: An exploratory case study Aleš Kresta, Michaela Sedláková

What makes stocks sensitive to investor sentiment: An analysis based on Google Trends Adeel Ali Qureshi

Financial inclusion, remittances and household consumption in sub-Saharan Africa: Evidence from the application of an endogenous threshold dynamic panel model Mahamat Ibrahim Ahmat-Tidjani

Economic growth in the European Union: Exploring the role of innovation and gender Vicente J. Coronel, Carmen Díaz-Roldán

Game-theory behaviour of large language models: The case of Keynesian beauty contests Siting Estee Lu

Editorial Board

Monika Banaszewska (Editor-in-Chief), Ivo Bischoff, Horst Brezinski, Gary L. Evans, Niels Hermes, Witold Jurek, Tadeusz Kowalski, Joanna Lizińska, Ida Musiałkowska, Paweł Niszczota, Michał Pilc (Deputy Editor-in-Chief), Katarzyna Schmidt-Jessa, Konrad Sobański

International Editorial Advisory Board

Edward I. Altman – NYU Stern School of Business Udo Broll – School of International Studies (ZIS). Technische Universität. Dresden Conrad Ciccotello – University of Denver, Denver Wojciech Florkowski – University of Georgia, Griffin Oded Galor – Brown University, Providence Binam Ghimire – Northumbria University, Newcastle upon Tyne Christopher J. Green – Loughborough University Eduard Hochreiter – The Vienna Institute for International Economic Studies Mark J. Holmes – University of Waikato, Hamilton Andreas Irmen – University of Luxembourg Bruce E. Kaufman – Georgia State University, Atlanta Robert Lensink – University of Groningen Steve Letza – The European Centre for Corporate Governance Robert McMaster – University of Glasgow Victor Murinde – SOAS University of London Hugh Scullion – National University of Ireland, Galway Yochanan Shachmurove – The City College, City University of New York Thomas Taylor - School of Business and Accountancy, Wake Forest University, Winston-Salem Linda Goncalves Veiga – University of Minho, Braga Thomas D. Willett – Claremont Graduate University and Claremont McKenna College Habte G. Woldu - School of Management, The University of Texas at Dallas

Thematic Editors

Economics: Monika Banaszewska, Ivo Bischoff, Horst Brezinski, Niels Hermes, Witold Jurek, Tadeusz Kowalski, Ida Musiałkowska, Michał Pilc, Konrad Sobański • Finance: Monika Banaszewska, Gary Evans, Witold Jurek, Joanna Lizińska, Paweł Niszczota, Katarzyna Schmidt-Jessa, Konrad Sobański • Statistics: Marcin Anholcer, Maciej Beręsewicz, Elżbieta Gołata

Language Editor: Robert Pagett

Paper based publication

© Copyright by Authors, Poznań 2025
 © Copyright for this edition by Poznań University of Economics and Business, Poznań 2025



This work is licensed under a Creative Commons Attribution 4.0 International License https://creativecommons.org/licenses/by/4.0

https://doi.org/10.18559/ebr.2025.2

ISSN 2392-1641 e-ISSN 2450-0097

POZNAŃ UNIVERSITY OF ECONOMICS AND BUSINESS PRESS ul. Powstańców Wielkopolskich 16, 61-895 Poznań, Poland phone +48 61 854 31 54, +48 61 854 31 55 https://wydawnictwo.ue.poznan.pl, e-mail: wydawnictwo@ue.poznan.pl postal address: al. Niepodległości 10, 61-875 Poznań, Poland

Printed and bound in Poland by: Poznań University of Economics and Business Print Shop

Circulation: 80 copies



Economics and Business Review

Volume 11 (2) 2025

CONTENTS

Editorial introduction	
Joanna Lizińska, Katarzyna Schmidt-Jessa, Konrad Sobański	3

ARTICLES

How initial price history influences expectation formation in multi-asset experimental markets: An exploratory case study	
Aleš Kresta, Michaela Sedláková	7
What makes stocks sensitive to investor sentiment: An analysis based on Google Trends Adeel Ali Qureshi	39
Financial inclusion, remittances and household consumption in sub-Saharan Africa: Evidence from the application of an endogenous threshold dynamic panel model Mahamat Ibrahim Ahmat-Tidjani	67
Economic growth in the European Union: Exploring the role of innovation and gender Vicente J. Coronel, Carmen Díaz-Roldán	91
Game-theory behaviour of large language models: The case of Keynesian beauty contests Siting Estee Lu	119
String Estee Education 1	

Game-theory behaviour of large language models: The case of Keynesian beauty contests



Abstract	Keywords
The growing adoption of large language models (LLMs) pre- sents potential for deeper understanding of human behav- iours within game theory frameworks. This paper examines strategic interactions among multiple types of LLM-based agents in a classical beauty contest game. LLM-based agents demonstrate varying depth of reasoning that fall within a range of level-0 to 1, which are lower than experimental results conducted with human subjects in previous studies. However, they do display a similar convergence pattern to- wards Nash Equilibrium choice in repeated settings. Through simulations that vary the group composition of agent types, I found that environments with a lower strategic uncertainty enhance convergence for LLM-based agents, and environ- ments with mixed strategic types accelerate convergence for all. Results with simulated agents not only convey in- sights into potential human behaviours in competitive set- tings, but also prove valuable for understanding strategic interactions among algorithms.	 large language models economic games strategic interactions
JEL codes: C63, C70, C90	
Article received 29 March 2025, accepted 9 June 2025.	

Suggested citation: Lu, S. E. (2025). Game-theory behaviour of large language models: The case of Keynesian beauty contests. *Economics and Business Review*, *11*(2), 119–148. https://doi.org/10.18559/ebr.2025.2.2182



This work is licensed under a Creative Commons Attribution 4.0 International License https://creativecommons.org/licenses/by/4.0

¹ School of Economics, University of Edinburgh – 30 Buccleuch Pl, Edinburgh EH8 9JT, United Kingdom, estee.I828@gmail.com, https://orcid.org/0009-0005-1212-5382.

Introduction

With the emergent line of research into large language models' (LLMs) capabilities, there are also growing discussions on the implications of LLMs for economic research and social sciences experiments, particularly in the field of game theory. One of this work's main objectives is to make a case for using LLMs as synthetic agents in economic games to shed light on potential strategic behaviours. Since LLMs are trained using human-generated data, observing interactions between them could be relatable to human subjects in experiments, and offer more insights than conventional simulation methods. As opposed to diving into more expensive human-based experiments straightaway, it is also relatively easy and cost-effective to test different setups before concentrating on designs that are worth pursuing.

Previous studies mainly focused on exploring two-player cooperative and non-cooperative games, and they often consist of a single LLM type (Akata et al., 2023; Horton, 2023; Phelps & Russell, 2023). While they provide interesting baselines for evaluating strategic behaviours, assuming agent homogeneity could make behaviour modelling more restrictive and does not leverage the potential of having multiple LLMs in the market. Furthermore, competitive games involve more strategic consideration in attempting to predict and outmanoeuvre opponents. Therefore, exploring such games could offer new insights into strategic interactions that are different from other games, providing novel and promising applications for LLMs. As a result, in this paper, I investigate a classical multi-player competitive game widely studied in economics: the beauty contest game. In this framework, agents' strategic levels and adaptive learning behaviours can be jointly evaluated. The methodology builds on top of a well-established line of research, thus providing a solid foundation for the approach adopted, complemented by the availability of human subject experimental results and broad economic applications to draw relation to.

In the first experiment involving multiple LLM types, I found that LLMbased agents manifest strategic levels between 0 and 1, evaluated using Nagel's (1995) level-k model, which are lower than experimental findings conducted with human subjects in previous literature. However, in repeated beauty contests with revelation of past information, most LLM-based agents show convergence towards the Nash equilibrium (NE) choice, mirroring that of human subjects.

Furthermore, I also explore simulations of beauty contest games in different game environments. Since opponent types could be important in influencing adaptive learning behaviour, I varied the proportion of agents with different strategic types within the group to analyse their impact on game outcomes. I found that when facing fixed-strategy opponents, LLM-based agents display faster convergence in low strategic uncertainty environments. When two types of LLM-based agents, one with higher strategic level than the other, are playing against each other, all agents display faster learning speed in such mixed environments than when they are playing against their own types. These results contribute to assessing LLMs with human-based metrics on strategic levels, thereby allowing for representation of heterogeneous human subjects with different LLM types. Potential strategic behaviours can also be explored via simulation of various set-ups, and postulating the possible implications.

On a broader view, given LLMs' capability, this work not only seeks to contribute to the growing literature on using LLM-based agents as a tool for social science research, and in simulating and deciphering human's strategic behaviours. I show that theories that were developed to explain and evaluate human behaviours can unequivocally help us to understand how this new era of computer algorithms would behave when competing against each other. With the growing integration of LLMs into daily life, where they can be used as surrogate agents to communicate and interact with one another, understanding how algorithms react to each other could have significant social impacts and real-world applications, particularly for competitive games, such as beauty contests.

The rest of the paper is structured as follows: Section 1 highlights the background. Section 2 explores the one-shot and repeated beauty contest games. Section 3 further investigates LLMs' adaptive learning behaviour via simulations of beauty contests with variation in group composition. Section 4 discusses the limitations and extensions, followed by conclusion.

1. Background

LLMs as a computational model of human behaviour. Since the training process of LLMs uses human-generated data and refinements based on direct human feedback, human reasoning process are baked into the algorithms (OpenAI, 2024; Ouyang et al., 2022). Therefore, it is proposed that LLMs can be perceived as an implicit computational model of human behaviour (Horton, 2023). I hereby streamline and differentiate between the two main aspects of how LLMs' human-like behaviour could apply to research for the social sciences community:

(a) Imitation of decision-making with known constraints. One approach is to use LLMs as synthetic agents with pre-specified profiles. The objective is to granulate the elements contributing to decision-making by testing outcomes given known constraints. This resembles agent-based modelling (ABM) (Hamill & Gilbert, 2015), where agents are pre-programmed to behave as we expect, and the outcome serves as a form of visualising and checking theoretical predictions. Applying this approach to the beauty contest games implies setting the strategic levels of the LLM-based agents *a priori* and examining their behaviours in comparison to theoretical predictions of agents with a certain strategic level.

(b) Mirroring human-like behaviours without known constraints. By abstracting away from putting restrictions on behaviours *a priori*, simulations conducted with LLM-based agents essentially offer a tool for computational experiments. In the context of beauty contests, this approach identifies the intrinsic strategic levels of the LLM-based agents, given a pre-specified game environment. By varying the experimental design, the behaviours of LLMbased agents could be used as pilots. The results can form conjectures of the possible outcomes if the experiments were conducted with human subjects.

In this paper, I focus on the second approach, which is more relevant to my objective of simulating potential strategic behaviours between LLM-based agents in a competitive setting. Furthermore, this method also accounts for the potential changes in strategic levels over time in repeated settings, which would not be identified if strategic levels are pre-fixed, as in (a).

LLMs as heterogeneous agents. Existing works (e.g., Akata et al., 2023; Horton, 2023) mainly explore the use of a single type of LLM to represent agents and do not fully leverage the potential of many different LLMs in the market. The presence of multiple LLMs could be used to model games with heterogeneous agents. There are many ways to define agent heterogeneity, one of which could be based on differences in the underlying training data. For instance, Anthropic's reward model training data primarily comes from crowd-sourcing feedback through Amazon Mechanical Turk, a platform often used for social sciences research; and OpenAl's models are mainly trained on used prompts (HuggingFace, 2022). LLMs could also comprise of different priors and come in varying sizes, leading to different performances in textbased generating ability. Therefore, each LLM can be perceived as representing a different type of agent. As a result, the LLMs used in this work comprise of models from different developers and of different sizes and architectures. However, while the above distinctions of types are intuitive and straightforward, they do not necessarily imply heterogeneity in strategic situations, which I seek to study. Therefore, I define LLM types by their corresponding strategic levels, determined through the one-shot beauty contest game using a measure ubiquitous to how we evaluate the strategic level of human subjects. This measure of agent heterogeneity also allows me to draw parallels between strategic behaviours displayed by LLM-based agents and different groups of human subjects. It also provides a flexible set-up where new models can be added and evaluated in a similar manner.

LLMs as complements to human participants. At the core of discussions surrounding the usefulness of LLMs in social sciences research is the question

of whether they can rise to the challenge of participating in social experiments in place of human subjects or as rational players. There are growing replications of social experiments and strategic games to investigate this. While it was found that LLM-based agents deviate away from game-theoretical predictions and may be far from rational, they inevitably demonstrate an ability to imitate human behaviours, making them human-like participants (Aher et al., 2023; Argyle et al., 2023; Dillion et al., 2023; Fan et al., 2023; Guo, 2023; Guo et al., 2024; Huijzer & Hill, 2023; Mei et al., 2024; Webb et al., 2023).

The main concern about using LLM-based agents is the opacity of their minds, which makes interpretations about their beliefs superficial (Dillion et al., 2023). Although the same argument applies to human minds, there exist many theories to describe human reasoning in strategic situations, but a lack of any equivalent to decipher the "thinking" process of AI algorithms. However, since LLMs are trained on human-generated data, which includes reasoning procedures, they could develop mechanisms similar to those of the human brain, thus theories applied to humans might also be applicable for explaining behaviours displayed by LLM-based agents (Kosinski, 2023). Furthermore, Strachan et al. (2024) measure LLMs' theory of mind ability and show that these could be on a par with or even outperform humans in terms of the ability to understand others' mental states, reflective of reasoning ability. This implied eliciting of reasoning from LLM-based agents could illuminate decision-making process undertaken by human subjects. However, despite this connection, given the opacity of both LLM-based agents and human subjects' internal reasoning processes, it remains important to treat simulated results with care, thus my work focuses more on revealed choices than the reasoning process. It does not aim to argue for replacing human subjects in experiments with LLM-based agents completely, but rather using them as complements to shed some light on potential strategic behaviours.

Choice of beauty contests. In this paper, I focus specifically on a beauty contest game, contributing to the study on multi-player competitive games with LLM-based agents. This set-up is desirable, as it encompasses both competitive nature and interactions between multiple, and possibly heterogeneous, agents, whose level of reasoning can be easily distinguished (Camerer et al., 2004; Nagel, 1995). The game can also be constructed with a single interior NE solution, even in repeated settings, obstructing away from the complication of analysing multiple equilibria. Furthermore, there are many applications of beauty contest games with substantial social value. For instance, the Keynesian Beauty Contest started off with a practical application to describe the stock market (Keynes, 1936; Nagel et al., 2017). With the market becoming more computerised, crypto trading bots emerge and function by executing pre-defined buying and selling strategies (Trality, 2024). The backbone of these automatic bots can be replaced in the future by LLMs that account for vast human data on trading behaviours, and one could instead focus on

choosing between different LLMs that behave as proxies for human traders. Therefore, understanding LLM interactions could better inform us about the potential social implications, such as in the trading market, and a beauty contest game is a good starting point.

2. Beauty contest games

In this section, I first explore the one-shot and repeated beauty contests involving multiple LLMs: ChatGLM2, ChatGLM3, Llama2, Baichuan2, Claude1, Claude2, PaLM, GPT3.5, GPT4. I will focus my analysis on determining the strategic levels associated with each LLM-based agent, and explore their learning patterns over time.

The results are based on experimental data adapted from Guo et al. (2024). However, unlike Guo et al. (2024), whose main objective was to evaluate LLMs' performance relative to rational players that select the NE choice, this work aims to analyse LLMs' behaviour as though they were human players.

General experimental design. Using a modified set-up following Nagel (1995), and an exemplary prompt, following Guo et al. (2024) (recited in Appendix A1):

Agents are asked to choose a number between 0 and \overline{c} , where \overline{c} is randomly generated from 0 to 1,000. The agent choosing the number closest to p, p = 2/3, of the average wins the game. A fixed prize of x is awarded to the winner. The prize is split amongst those who tie.

In a repeated beauty contest game, the same game is played for 6 periods, and agents are given historical information up to 3 past periods. These include choices made by all agents, the average of these choices, 2/3 of the average, and past winners. The choice of revealing up to 3 past periods is due to token restrictions to control computation intensity. As a result, this set-up can be perceived as one with partial feedback or an exogenous forgetting parameter.

Data collection. The experiments are conducted with API calls of different LLMs, providing a collection of independent observations that allows for a robust measure of strategic level for each LLM type. In repeated settings, the information availability can be explicitly controlled through prompts that reveal histories perfectly or selectively to LLMs (Bauer et al., 2023). While the stochasticity of model responses is dependent on the temperature selected, Chen et al. (2023) show that strategic or choice consistency is less influenced by temperature, which depends more on the underlying reasoning process. Therefore, this work does not explore changes in responses given variations in temperature.

Analysis Focus. The two main concepts central to my analysis are:

- Determination of strategic levels. Following Nagel (1995), an agent is of strategic degree *n* if it chooses a number $r(2/3)^n$, where *r* is defined to be the reference point, characterised by naive player or a point of salience in heuristics. In one-shot games and in period 1 of repeated games, this reference point is assumed to be the mean of the range of numbers in the action space (ie. $r = \overline{c}/2$).
- Convergence. In repeated games, changes in choices are tracked to determine if there is convergence to the unique NE of 0. The convergence rate is computed as $c_t = -(a_{t+1} a_t)/a_t$, where $a_{t+1} \le a_t$, a_t being the action/number chosen in period t. Changes in strategic levels are found by re-adjusting the reference point to the mean of the previous period choices.

2.1. One-shot game

150 sessions of one-shot beauty contest were conducted with 9 agents represented by different LLMs. In classical beauty contests, \overline{c} is often fixed at 100, and as a result, all choices between (66.66, 100] are weakly dominated by 66.66, and those above 44.44 are weakly dominated by 44.44, etc. Via iterative elimination of weakly dominated strategies, the number of steps taken determines agents' strategic levels. Otherwise, going by the level-k model with a focal point set at the mean of the number range, 50, level-0 would choose 50, and level-1 responds by choosing 33.33, etc. The unique interior *NE* solution of the game is 0. In this modified set-up with a randomly generated upper bound for each game, the steps of assessing the strategic levels are unaffected. For example, using the level-k model, level-0 would simply choose the focal point, $\overline{c}/2$, and level-1 would respond by choosing 2/3 $\cdot \overline{c}/2$.

Choices. Figure 1 shows that the normalised choices made by LLM-based agents are concentrated at 50 for ChatGLM3, Baichuan2, Claude1, PaLM. As per level-*k* model, they are level-0 players. Llama2 records fairly dispersed and randomised choices, and thus can be perceived as level-0 as well. Claude2 shows a spike around 33, indicating likelihood of level-1 thinking. There is also high choice frequency around 66, which could be rationalised as step-1 of iterated elimination of dominated strategies (Mauersberger & Nagel, 2018). For GPT3.5, most of the choices are concentrated around 33, stipulating level-1 reasoning. While there are some other spikes at 50 and 66, those are of much lower frequency. GPT4 displays the highest spike in choices around 44, implying step-2 depth of reasoning by iterated elimination of dominated strategies. A lower spike is also observed around 33, corresponding to level-1 thinking in the level-*k* model. This could suggest that GPT4 has a level in between 1



Figure 1. Many LLM-based agents choose 50 with higher frequency

and 2. No data is observed for ChatGLM2, indicating it is unable to complete the games and produce comprehensible output given the instructions.

Nagel (1995) and Bosch-Domenech et al. (2002) have conducted beauty contest games with different human populations, such as students (mean = 36.73, median = 33), theorist (mean = 17.15, median = 15*), newspaper readers (mean = 23.08, median = 22*), etc.² In their studies, human subjects show a strong deviation away from game-theoretic prediction, and display on average iteration steps 1 and 2 evaluated by the level-k model. Compared to them, LLM-based agents choose slightly higher numbers, as shown in Table 1, which corresponds to an average strategic level between 0 to 1. This could be due to differences in human subjects involved in the experiments and the underlying data used to train the LLMs. Moreover, since LLM-based agents could display different strategic levels, their behaviour could be representative of different subsets of the population.

Models	Chat GLM3	Chat GLM2	Llama2	Bai- chuan2	Claude2	Claude1	PaLM	GPT3.5	GPT4
Average	52.029	N/A	59.519	51.158	41.609	47.696	49.976	38.912	41.072
Median	51.724	N/A	62.685	50.0	33.333	49.313	50.0	33.333	44.442

Table 1. Average and median choice of LLM-based agents across 150 Sessions

Source: own work.

For human subjects, when given an identical game set-up, it is possible that they might employ different strategies (Costa-Gomes & Weizsäcker, 2008; Devetag et al., 2016). The same could apply to LLM-based agents. Therefore, by fixing the game parameters and instructions, it is possible to analyse how varied agents' choices might be.

Figure 2 shows that within the 150 sessions, for the same upper-bound value, \overline{c} , Claude2, GPT3.5 and GPT4 displayed more variability in choices than other models. This is similar to human players, where choices might not be static even when the game parameters and instructions are the same, LLM agents' behaviour also encompass this aspect to some extent. While Bauer et al. (2023) indicate that running multiple sessions could already accommodate the stochastic nature of LLM responses, my method of using average choices based on both identical and different upper bounds could render a more robust and consistent measure of strategic levels for each model.

Strategic levels. Following the level-k model to compute for the strategic levels, \overline{n} , the reference point, r, is defined to be the choice of a non-strategic agent, which is assumed to be the mean of the number range, pertaining

 $^{^{\}rm 2}$ The median with * are guesstimated based on the figures in Nagel (1995) and Bosch-Domenech et al. (2002).



Figure 2. Some LLMs display variability in chosen number given the same upper bound

to insufficient reasoning (Mauersberger & Nagel, 2018). However, this focal point can be disputable. In my set-up, the varied upper bounds may also be the focal points rather than taking the extra step of computing for the mean. In Figure 3, I show that the average strategic levels are between 0 and 1 given the reference point $r = \overline{c}/2$, and between 1 to 2.5 when it is $r = \overline{c}$. However, for consistency with the existing literature on beauty contests, in the following sections, I evaluate the results using the conventional focal point of $\overline{c}/2$.

Comparing across LLM-based agents, in Figure 3a, the strategic levels are relatively high for ChatGLM3, Claude2, GPT3.5 and GPT4. Surprisingly, GPT4 has a slightly lower strategic level than GPT3.5, even though it is often presumed to be a stronger model. It may be possible that its lower depth of reasoning is due to it being trained on more data, thereby encompassing a higher possibility of noisy strategies, leading to a higher average chosen number.

Figure 3 also shows variability in strategic levels, which could again indicate some degree of choice inconsistency that is similar to human subjects. While this highlights the plausibility of exploring agent heterogeneity on another dimension of variability in strategic levels, this work follows a conventional analysis approach in beauty contests and focuses on average strategic levels.

Payoff. Figure 4 demonstrates that Claude2, GPT3.5 and GPT4 have relatively higher average payoffs than the others, of which GPT3.5 has the highest average payoffs compared to the other models. Associating the results with strategic levels, LLM-based agents with higher average strategic levels can often obtain higher average payoffs, except for ChatGLM3. This could be







Figure 3. Average strategic levels of LLM-based agents with reference point $r = \overline{c}/2$ (in 3a) fall between 0 to 1, and for $r = \overline{c}$ (in 3b), they are between 1 to 2.5



Figure 4. Average payoffs are higher for Claude2, GPT3.5 and GPT4

Source: own work.



Figure 5. Most LLM-based agents display convergence in average chosen number

due to high variability in the strategic level of the ChatGLM3-based agent, thus adversely influencing its average gain.

2.2. Repeated games

Following the repeated set-up highlighted in the general experimental design, 30 sessions of repeated beauty contests were conducted.

In Figure 5, most LLM-based agents show convergence in actions, particularly for Claude1, Claude2, GPT3.5, and GPT4, which are models of higher strategic levels, determined by the one-shot games. Their chosen numbers are approximately 0 in period 6, indicative of them learning to play *NE* choice across time.

Evolution of strategic levels. Figure 6a shows the changes in strategic level across time for each LLM-based agent, averaged across sessions. While the strategic levels evolve over time, the range of change is narrow. On average, they stay within the bound of 0 and 1.4. Most LLM-based agents display increasing depth of reasoning, especially Claude2, GPT3.5 and GPT4. An interesting observation is that while GPT3.5 has a higher strategic level than GPT4 in one-shot games, in repeated settings, GPT4's average strategic level surpasses that of GPT3.5 from periods 2 onwards, implying that it could be more adept at revising its beliefs about opponents over time given past information. The abnormality in Figure 6a comes from ChatGLM3 and Llama2, the first shows a decrease in the average strategic level, indicating a lack of

ability to respond to historical information and adjust behaviour accordingly; the second displays naive, random selection throughout the periods, and on average, it fails to complete the game beyond period 4.

Payoff evolution. Figure 6b shows the transition of average payoff over time. GPT3.5 outperforms the other LLM-based agents in all periods; Claude2 and GPT4 also perform relatively well and they are more or less comparable; the rest of the LLM-based agents do not obtain as high an average payoff, but most of them display growth over time. Coupled with Figure 5, which shows convergence in average choice towards *NE*, the increasing payoffs could be an indication of learning about the optimal action to take in order to win the game.

In this section, the purpose of evaluating the one-shot games is to determine the strategic levels of LLM-based agents. The computation method is the one conventionally used in human subject experiments, and thus allows par-



Figure 6. Average strategic levels (6a) and average payoffs (6b) across 30 sessions for 6 periods are highest for GPT4 and GPT3.5 respectively

Source: own work.

allels to be drawn between the results. Experiments with LLM-based agents resemble those conducted with human subjects: they both show strong deviation away from game-theoretic prediction, and agents tend to display low levels of reasoning. However, the distinction is that LLM-based agents display an even lower level of reasoning as compared to human subjects. Furthermore, the repeated setting sheds light on how simulated agents could learn over time. In a similar way to human subjects, LLM-based agents do not display iteration steps that go over 2 within the span of the games, but they do seem to learn from historical information and show convergence towards NE choice.

3. Simulation of adaptive learning behaviour with variation in group composition

Following on from above, in this section, I explore LLM-based agents' learning patterns further by analysing how variations in group composition could affect their behaviours. These results can also be perceived as computational experiments conducted with synthetic agents, which may illuminate human behaviour in similar set-ups and would be useful as insights for experimental pilots.

Based on the strategic levels determined, I choose two LLMs with different strategic levels, GPT3.5 and PaLM. GPT3.5 has a strategic level of approximately 1 and PaLM has level-0, representing a higher (H) and lower (L) intelligence agent type, respectively, where intelligence is interpreted loosely as a metonym for strategic level. I will construct groups of heterogeneous agents using these two types of LLM-based agents.

Set-up. Games are played among 10 agents, who are asked to choose a number between [0, 100]. The same group plays for 5 periods with full historical information disclosure (i.e. choices made by all agents, average of these choices, 2/3 of the average, and past winners). The winner is the agent whose number is the closest to 2/3 times the average of all chosen numbers. In each period, the winner receives a fixed prize of x. In the case of a tie, the prize is split amongst those who tie, and all other players receive 0.

3.1. Partial static environment: LLM vs. static algorithm

In this environment, LLM-based agents are asked to play against fixed-strategy players, whose actions are hard-coded to be 0. There are 3 treatments: (1) 1 LLM + 9 Hard-coded Agents (Low strategic uncertainty); (2) 5 LLMs + 5 Hard-coded Agents (Mixed strategic uncertainty); (3) 9 LLMs + 1 Hard-coded Agents (High strategic uncertainty). These treatments allow analysis of agents' behaviour amidst different levels of strategic uncertainty. Across different treatments, the proportion of fixed strategy players and LLM-based agents change, but the group size remains the same. LLM-based agents are also told that some of their opponents are playing a fixed strategy of 0. An exemplary prompt is shown in Appendix A2.

For both types of LLM-based agents, there is convergence in choices to 0 in general, as shown in Figure 7 and 8. This learning pattern exhibits either refinement of beliefs about opponents' strategies or progression in their own depth of strategic thinking when given historical information. The pace is slower as strategic uncertainty grows, where the proportion of LLM-based agents becomes larger relative to fixed-strategy players.

Comparing the high (H) and low (L) types, all H agents chose the same number over time in Treatment 2 and 3, where there are multiple LLM-based agents. Therefore, they are shown in Figure 7 as representative agents. However, not all L agents choose the same number in those treatments, as shown by multiple graphs in Figure 8b and 8c, which indicates that some Lagents may choose different numbers. This demonstrates that when strategic uncertainty is high, L displays larger variability in choices and there might not be any convergence at all.

Furthermore, L types also behave less "cautiously" in the sense that they could converge to 0 in period 2 straightaway when strategic uncertainty is relatively low, whereas convergence to 0 takes a gradual process for H. This could indicate that H goes from less sophisticated strategies to more refined choices through iterative learning and adaptation, and there is a lack of such system-



Figure 7. Choices of higher intelligence LLM-based agents playing against fixed strategy opponents display gradual convergence in all treatments

Source: own work.



Figure 8. Choices of lower intelligence LLM-based agents playing against fixed strategy opponents for Treatment 1 (8a), 2 (8b), and 3 (8c) may display abrupt adjustment or lack of convergence

atic adjustments in choices for *L*, which could suggest that they are relying more on intuitive guesses than successive elimination of less likely options.

Evolution of strategic levels. When evaluating the transition in strategic levels across periods, H shows the transition from 0 to 1, and most of L-type agents stay at level-0, with some fluctuations between 0 and 1 when strategic uncertainty is high.

Payoffs. The payoffs are in favour of LLM-based agents rather than fixed-strategy players when strategic uncertainty is relatively high. H could gain better payoffs in a low and high strategic uncertainty environment as compared to a mixed strategic uncertainty environment, where they receive a flat payoff of 0 throughout the periods. Comparing the types, it is interesting to note that payoffs achieved by L in all settings may be comparable or even higher than that of H, although the variations are also larger. This indicates that a higher strategic level does not necessarily imply higher payoffs when competing against fixed strategy opponents. These results not only signify the potential game play if human subjects are playing against opponents that naively adopt a fixed strategy of 0, but could also illustrate a possible outcome if they are going against static computer algorithms executing a fixed NE strategy (see Appendix A3.2, Figure A1 & A2).

Application. One example of beauty contest applications is the Bertrand competition model (Mauersberger & Nagel, 2018). LLM-based and fixed strategy agents can be perceived as firms adopting different pricing strategies, with the objective being to win over the market and maximise their profits. Fixed strategy firms could be perceived as playing the equilibrium action by setting the price equals to marginal cost, while LLM-based firms could be more dynamic and adjust their prices in each period.

In terms of payoffs, if there exists some rigidity in the short run, such as production capacity constraints for the firms or limited response time for the consumers, then those who set higher prices would be able to gain higher profits. In the long run, however, all factor inputs are flexible and consumers will not purchase from a firm that sells a homogeneous product at a higher price than the equilibrium. As a result, H-type firms could often achieve better outcomes than L-type ones in the short run, where they can earn a positive profit by converging gradually. Even in the long run, the larger variance in pricing strategies for the L type as compared to H could result in them failing to converge to the NE, or in them displaying higher volatility in pricing, both of which could adversely impact their profits.

If firms outsource their pricing strategies to automated algorithms, this simulation could also be interpreted as competition between different algorithms. While automated pricing has been widely discussed in literature, those represented by LLMs that could respond to changes in rivals' strategies by adjusting their own ones could spark fresh perspectives (Brown & MacKay, 2023; Chen et al., 2016).

3.2. Dynamic environment: LLM vs. LLM

In this setting, LLM-based agents are playing against each other (similarly, GPT3.5 is denoted as H, and PaLM as L). There are 5 treatments: (1) 10 H LLMs; (2) 9 H LLMs + 1 L LLM; (3) 5 H LLMs + 5 L LLMs; (4) 1 H LLM + 9 L LLMs; (5) 10 L LLMs. I use the original prompt as shown in Appendix A1.



Figure 9. Transition of chosen number given variation in group composition for LLM vs. LLM-based agents for different environments, including Pure High Intelligence (9a), Highly Intelligent (9b), Mixed Intelligent (9c), Less Intelligent (9d), Pure Low Intelligence (9e)

Source: own work.

In Figure 9, set-up 1 (Figure 9a) and set-up 5 (Figure 9e) depict pure intelligence environments. While H agents show adjustment in their choices to lower numbers, the L agents persistently choose around 50. Whereas in setup 2 to 4, both H and L agents show convergence to lower numbers. The main difference is that the gap between the numbers chosen by H and L is smaller when there is higher proportion of L agents in the group. This result shows that L agents fail to adapt their strategies in the pure environment despite given historical information, but when placed in environments with mixed types, these could instigate faster learning. This observation applies for both H and L types, particularly when there is a higher proportion of H agents. This puts forth the strong statement that adding a single H agent could very well speed up learning. Convergence of choices and evolution of strategic levels. Figure 10 shows for set-up 1 and 5, the convergence rates for choices are low and approximately flat. In the mixed environments, the convergence speed fluctuates but could be higher than the pure environments. For instance, most of the convergence rates in set-up 2 to 4 lie above the lines for set-up 1 and 5. The higher the proportion of H, the higher the convergence rates. When computing for variations in strategic levels across time, all set-ups except for 5, where L agents do not display any apparent evidence of learning, show changes in strategic levels. In set-up 3, in particular, H could reach a strategic level greater than 1, which implies that having a highly mixed environment could also stimulate considerable growth in terms of depth of reasoning for some agents. A possible conjecture for this could be that the strategic landscape is more complex in a highly mixed environment: agents cannot simply default to strategies assuming similar reasoning process from all agents, and this may induce increasing depth of reasoning.



Figure 10. Average convergence rates are low and approximately flat for pure type environments

Source: own work.

Payoffs. The maximum possible payoffs that can be achieved in the mixed environment is either comparable or could be higher than that of pure environments. Since this is a competitive game, a higher gain for some agents also means higher losses for some, thus the variability in payoff outcomes, even for the same agent type, can also larger. While *L* agents usually obtain positive payoffs at the beginning of the game for choosing 50, which is closer to 2/3 of the average, this head-start is soon eroded if the group contains any *H* agents, who learn to react to this information rapidly. Therefore, *L* agents are less likely to win across periods. Furthermore, the degree of heterogene-

ity also matters. H agents could obtain higher average payoffs at the expense of L agents when $L \ge 50\%$, and L agents are better off if there are less H (see Appendix A3.2, Figure A3).

Application. The simulation results could assist in informing policies. A potential application is the streaming system in schools, where students are allocated into different classes based on their grades to facilitate better learning (Ireson & Hallam, 1999; Liem et al., 2013). Let us suppose students are classified into high and low types in terms of ability: my findings provide an argument for a mixed learning environment, where the low types would learn faster when integrated into a class with larger proportion of higher ability peers; even for high types, their learning rate could be slightly improved.

Furthermore, the results also make a case for the usefulness of sustaining a variety of LLMs, including weaker models. Even though they do not learn when competing against each other, they could learn when placed in the presence of stronger LLMs. Stronger LLMs could also benefit from playing against a small proportion of weaker LLMs, as shown by higher learning rates, and they could also achieve better average payoffs when playing against larger proportion of weaker LLMs. This suggests the value of continual investments in LLMs of differing strength.

Reasoning elicitation. It is recognised that drawing direct relations between LLM-based agents and humans in terms of internal reasoning process may be speculative and overextending parallels, therefore analysing observed actions takes precedence in this paper. However, with the growing body of literature that highlights LLMs exhibit human-like reasoning (e.g., Kosinski, 2023; Strachan et al., 2024), eliciting reasoning in computational experiments may serve as an avenue to gain some perspective on agents' rationales for making certain choices and how they might learn.

In all set-ups, LLM-based agents were prompted at the beginning of period 1 to state their understanding of the game, and for each subsequent periods, they are asked to restate the goals. This step is essential to mitigate the potential of them not comprehending the game, in which case, all LLM-based agents are able to correctly recite the game rules.

The agents were also asked to give a statement of reasoning in support of their choices. In period 1, both H and L agents make choices based on their belief of a popular number, which is often the mean of the range. In subsequent periods, L agents appear to learn by either adjusting the reference point, and make selections that still comply with a strategic level of 0, or via imitation by following the winner's past choice. They may also not learn at all, and continue to select a number that they believe to be the popular choice. As for the H agents, they can learn by (1) anchoring their guesses to two-thirds of the past period's average; (2) imitating the winner's strategy; (3) making adjustments based on past period payoffs; and also (4) pattern recognition. Agents may place different reliance on distinct pieces of historical informa-

tion when making their choices, and multiple types of learning could come into play. This diversity in learning mechanisms could lead to higher speed of changes in average choices, and in turn translate into a higher strategic level.

4. Limitations and extensions

Much like experiments with human subjects, LLM-based agents could also be sensitive to variations in game design, feedback, and instructions. This work only explored a small number of set-ups and for a particular competitive game, which can limit the scope. However, it serves the main purpose of proposing the potential of LLMs as a valuable tool for social sciences research, and beauty contests being a game of substantial impact in economics research provide an excellent foundation for this line of work. The simulation results not only shed light on potential strategic behaviours given variations in set-ups, they also illuminate outcomes when algorithms are interacting with one another.

Some of the possible extensions would be to include:

Variations in game design and feedback. While I focused on p = 2/3, p can be varied to 1/2 or 4/3 to replicate human subject experiments, in which case, equilibrium multiplicity could arise, allowing for analysis on equilibrium selection (Nagel, 1995). In addition, the same set-ups can be implemented but with variations in terms of which piece(s) of historical information to reveal.

Objectives. Humans are sensitive to problem framing and phrasing of survey questions. Similarly, LLMs' decisions could be influenced by the formatting of prompts as well (Kalton & Schuman, 1982; Sclar et al., 2023; Tversky & Kahneman, 1981). This work explores how agents behave when the objectives are set to be winning the game and followed by maximising their payoffs, but in most economic models, the primary focus is usually on maximising utilities and then winning. In this competitive game, the winning strategy is also one that gives the best payoff, thus changing the sequence of objectives is unlikely to result in drastic differences in game outcomes, but could serve as a sanity check.

Prompt language. In Guo et al. (2024), the prompt language was changed to Mandarin Chinese in the multi-LLM-based agents setting. It was found that PaLM is unable to complete the games, indicating the potential difficulty in comprehending the instructions when they are given in another language. As for GPT3.5, it can complete the game in a Chinese setting but the choices are more clustered. The variance in strategies observed as compared to the English setting may reflect differences in strategic behaviours among different language users that the models are trained on, or it could stem from a significantly smaller availability of human-generated data in another language.

While the current work focuses on an English setting, future work could involve replicating the set-ups in other prompt languages to model heterogeneous populations in other dimensions. Nonetheless, this result also underscores the scarcity of literature on comparing experimental outcomes across human subjects from different language backgrounds, which could have important implications if the game is applied to different cultural and linguistic contexts.

Human–machine interactions. Previously, experimental designs involving computers usually comprised of pre-defined algorithms, and humans were found to display a higher degree of strategic reasoning when competing against fellow human opponents as opposed to computer algorithms (Coricelli & Nagel, 2009). Human vs. LLMs could offer a fresh form of human-machine interactions, as LLM-based agents could respond dynamically and switch their strategies based on historical information, thereby contributing to greater strategic uncertainty and complexity. Given that LLMs display some degree of learning abilities, they could also learn from playing with human subjects, making the interactions more intriguing to explore.

Future validity. Another important question would be the future validity of the results proposed by this paper. Here, the measures of strategic levels are robust to the changing game parameters, such as the upper bound of the choice range, which could serve as a form of sensitivity test and make the results more replicable under the same conditions. Apart from this, there is growing interest in exploring whether prompting LLMs with questions could make them more strategically sophisticated in the future, and therefore the results cannot be replicated. This work shows that within a given session, models converge towards NE choice if they gain exposure to past play information, which is indicative of their learning ability over time, offering the possibility of individuals training their own algorithms to better fit their preferences in different contexts and LLMs becoming more sophisticated in the future. However, since the experiments are conducted with effectively stagnant LLM versions, and the information provided to the LLMs during the experimental sessions is controlled, this allows the validity and replicability of the results under the same set-ups. If future versions of LLMs incorporate the questions asked by the individuals into their training, then new models could be relatively more sophisticated or on the contrary, less sophisticated due to the incorporation of more noisy data. This could give rise to more serious concerns over whom are the models aligning to, which is an open question for future exploration.

Conclusions

The contribution of this work is threefold. Firstly, it serves as part of the literature that seeks to make a case for integrating LLMs as tools for social sciences research. It then proposes adopting LLM-based agents in multi-player competitive games and explores the beauty contest game in particular, which has a wide range of economic applications. Drawing parallels to human subjects, LLM-based agents were found to have strategic levels in between 0 and 1, which is slightly lower than human subjects. Most of them also exhibit learning from historical information, showing convergence to the NE choice at varying rates, demonstrating either a revision in "beliefs" about their opponents, or increasing depth of reasoning. Similar to human subjects, though strategic levels evolve over time, the increase is minimal. Furthermore, to better understand strategic interactions and learning in varying environments, I simulate game play between LLM-based agents and fixed-strategy opponents, as well as among LLM-based agents. By varying the proportion of agent types in each group, I found LLM-based agents converge slower to 0 as the proportion of fixed-strategy agents decreases, demonstrating the impact of increased strategic uncertainty. Agents were also found to learn faster when placed in environments with mixed strategic types than environments comprised of a single strategic type. This postulates the potential for stimulating faster learning, particularly among less intelligent agents, by introducing heterogeneity into the groups. Last but not least, this work not only demonstrates potential strategic behaviours among human subjects in various set-ups through simulations, it also offers some insights into how different algorithms could behave when interacting with one another, showing potential outcomes if algorithms were to act as proxies for humans and applied to competitive situations.

There are many possible extensions and great potential for LLMs to be employed as toolkits for social sciences research in interpreting and deciphering human behaviour, which remains a relatively new subject area. The reverse is true as well, theories and experimental results from decades of learning about human decision-making can be similarly used to better understand machine behaviours and improve their performance.

Appendix

A1. Original Exemplary Prompt

One-shot games with multi-LLM-based agents

[System instruction.] I want you to act as a clever game player and perform as rationally as possible.

[User instruction.] You are playing a game and the rules are as follows:

- it consists of {number of players} players, including you,
- this is a one-round game,
- you need to choose a real number between 0 and {upper bound}, 0 and {upper bound} inclusive,
- everyone does not know how others would make choices beforehand,
- the winner is the person whose number is the closest to the two thirds of the average of all chosen numbers,
- the winner gets a fixed prize, in case of a tie the prize is split amongst those who tie,
- your goal is primarily to maximise the possibility of getting the prize and secondly to maximise your prize.

Subsequent prompt for historical information

[User instruction.]

The game of the same config has been held for {number of runs} run(s), and the historical choices of everyone are shown below (your id is {ID of the agent}: {historical information including (1) period index, (2) choices made by all agents, (3) average of the choices; (4) 2/3 of the average; (5) winner id.}

 Everyone can optimise his/her answer with the history to play in a new run in order to achieve goals.

(Return to Section 3.2).

A2. New exemplary prompt

Opponents playing fixed strategy of 0

[System instruction.] I want you to act as a clever game player and perform as rationally as possible.

[User instruction.] You are playing a game and the rules are as follows:

- it consists of {number of players} players, including you,
- this is a one-round game,
- you need to choose a real number between 0 and {upper bound}, 0 and {upper bound} inclusive,
- everyone does not know how others would make choices beforehand,
- the winner is the person whose number is the closest to the two thirds of the average of all chosen numbers,
- the winner gets a fixed prize, in case of a tie the prize is split amongst those who tie,
- your goal is primarily to maximise the possibility of getting the prize and secondly to maximise your prize,
- some of your opponents will be playing a fixed strategy of 0 and all others are behaving as rationally as possible.

Follow-up for each period.

Please just strictly output a JSON string, which has following keys:

- understanding: str, your brief understanding of the game,
- popular answer: float, the number which you think other players are most likely to choose,
- answer: float, the number which you would like to choose,
- reason: str, the brief reason why you give the popular answer and the answer that way.

Subsequent prompt (after period 1).

- The game of the same config has been held for {number of runs} run(s), and the historical choices of everyone are shown below (your id is {ID of the agent}: {historical information including (1) period index, (2) choices made by all agents, (3) average of the choices; (4) 2/3 of the average; (5) winner id}.
- Everyone can optimise his/her answer with the history to play in a new run in order to achieve goals.

(Return to Section 3.1).

A3. Additional details

For variations in group composition, I show below the payoff transition when playing against fixed strategy opponents:



Figure A1. Transition of payoffs for high type LLM-based agent(s) vs. fixed--strategy opponents in environments with Low Strategic Uncertainty (A1a), Mixed Strategic Uncertainty (A1b), and High Strategic Uncertainty (A1c)

Source: own work.



Figure A2. Transition of payoffs for low type LLM-based agent(s) vs. fixedstrategy opponents in environments with Low Strategic Uncertainty (A2a), Mixed Strategic Uncertainty (A2b), and High Strategic Uncertainty (A2c)

Source: own work.

I show below the payoff transition when playing with LLM-based opponents: (Return to Section 3).



Figure A3. Transition of payoffs given variation in group composition for LLM vs. LLM-based agents in environments with Pure High Intelligent (A3a), Highly Intelligent (A3b), Mixed Intelligent (A3c), Less Intelligent (A3d), and Pure Low Intelligent (A3e)

References

- Aher, G. V., Arriaga, R. I., & Kalai, A. T. (2023). Using large language models to simulate multiple humans and replicate human subject studies. *Proceedings of Machine Learning*, 202, 337–371. https://proceedings.mlr.press/v202/aher23a.html
- Akata, E., Schulz, L., Coda-Forno, J., Oh, S. J., Bethge, M., & Schulz, E. (2023). *Playing repeated games with large language models*. https://doi.org/10.48550/ arXiv.2305.16867
- Argyle, L. P., Busby, E. C., Fulda, N., Gubler, J. R., Rytting, C., & Wingate, D. (2023). Out of one, many: Using language models to simulate human samples. *Political Analysis*, *31*(3), 337–351. https://doi.org/10.1017/pan.2023.2
- Bauer, K., Liebich, L., Hinz, O., & Kosfeld, M. (2023). *Decoding GPT's hidden 'rationality'* of cooperation. SAFE Working Paper, 401. https://doi.org/10.2139/ssrn.4576036
- Bosch-Domenech, A., Montalvo, J. G., Nagel, R., & Satorra, A. (2002). One, two,(three), infinity,...: Newspaper and lab beauty-contest experiments. *American Economic Review*, 92(5), 1687–1701. https://doi.org/10.1257/000282802762024737
- Brown, Z. Y., & MacKay, A. (2023). Competition in pricing algorithms. American Economic Journal: Microeconomics, 15(2), 109–156. https://doi.org/10.1257/ mic.20210158
- Camerer, C. F., Ho, T. H., & Chong, J. K. (2004). A cognitive hierarchy model of games. *The Quarterly Journal of Economics*, 119(3), 861–898. https://doi. org/10.1162/0033553041502225
- Chen, L., Mislove, A., & Wilson, C. (2016). *An empirical analysis of algorithmic pricing on Amazon marketplace*. Proceedings of the 25th International Conference on World Wide Web, 1339–1349. https://doi.org/10.1145/2872427.2883089
- Chen, Y., Liu, T. X., Shan, Y., & Zhong, S. (2023). The emergence of economic rationality of GPT. *Proceedings of the National Academy of Sciences*, *120*(51), e2316205120. https://doi.org/10.1073/pnas.2316205120
- Coricelli, G., & Nagel, R. (2009). Neural correlates of depth of strategic reasoning in medial prefrontal cortex. *Proceedings of the National Academy of Sciences*, 106(23), 9163–9168. https://doi.org/10.1073/pnas.0807721106
- Costa-Gomes, M. A., & Weizsäcker, G. (2008). Stated beliefs and play in normal-form games. *The Review of Economic Studies*, *75*(3), 729–762. https://doi.org/10.1111/j.1467-937X.2008.00498.x
- Devetag, G., Di Guida, S., & Polonio, L. (2016). An eye-tracking study of feature-based choice in one-shot games. *Experimental Economics*, *19*(1), 177–201. https://doi. org/10.1007/s10683-015-9432-5
- Dillion, D., Tandon, N., Gu, Y., & Gray, K. (2023). Can ai language models replace human participants? *Trends in Cognitive Sciences*, *27*(7), 597–600. https://doi. org/10.1016/j.tics.2023.04.008
- Fan, C., Chen, J., Jin, Y., & He, H. (2023). Can large language models serve as rational players in game theory? A systematic analysis. https://doi.org/10.48550/ arXiv.2312.05488.
- Guo, F. (2023). GPT in game theory experiments. https://doi.org/10.48550/ arXiv.2305.05516.

- Guo, S., Bu, H., Wang, H., Ren, Y., Sui, D., Shang, Y., & Lu, S. (2024). *Economics arena for large language models*. https://doi.org/10.48550/arXiv.2401.01735.
- Hamill, L., & Gilbert, N. (2015). *Agent-based modelling in economics*. John Wiley & Sons.
- Horton, J. J. (2023). Large language models as simulated economic agents: What can we learn from homo silicus? NBER Working Paper, 31122.. https://doi.org/10.3386/ w31122
- HuggingFace. (2022). Illustrating reinforcement learning from human feedback (RLHF). https://huggingface.co/blog/rlhf
- Huijzer, R., & Hill, Y. (2023, January 31). Large language models show human behavior. https://doi.org/10.31234/osf.io/munc9
- Ireson, J., & Hallam, S. (1999). Raising standards: Is ability grouping the answer? Oxford Review of Education, 25(3), 343–358. https://doi.org/10.1080/030549899104026
- Kalton, G., & Schuman, H. (1982). The effect of the question on survey responses: A review. Journal of the Royal Statistical Society Series A, 145(1), 42–73. https:// doi.org/10.2307/2981421
- Keynes, J. M. (1936). The general theory of interest, employment and money. Macmillan.
- Kosinski, M. (2023). Theory of mind may have spontaneously emerged in large language models. https://www.gsb.stanford.edu/faculty-research/working-papers/ theory-mind-may-have-spontaneously-emerged-large-language-models
- Liem, G. A. D., Marsh, H. W., Martin, A. J., McInerney, D. M., & Yeung, A. S. (2013). The big-fish-little-pond effect and a national policy of within-school ability streaming: Alternative frames of reference. *American Educational Research Journal*, 50(2), 326–370. https://doi.org/10.3102/0002831212464511
- Mauersberger, F., & Nagel, R. (2018). Levels of reasoning in Keynesian beauty contests: A generative framework. In C. Hommes & B. LeBaron (Eds.), Handbook of computational economics (vol. 4, pp. 541–634). Elsevier. https://doi.org/10.1016/ bs.hescom.2018.05.002
- Mei, Q., Xie, Y., Yuan, W., & Jackson, M. O. (2024). A Turing test of whether ai chatbots are behaviorally similar to humans. *Proceedings of the National Academy of Sciences*, 121(9), e2313925121. https://doi.org/10.1073/pnas.2313925121
- Nagel, R. (1995). Unraveling in guessing games: An experimental study. *The American Economic Review*, *85*(5), 1313–1326. https://www.jstor.org/stable/2950991
- Nagel, R., Bühren, C., & Frank, B. (2017). Inspired and inspiring: Hervé moulin and the discovery of the beauty contest game. *Mathematical Social Sciences*, *90*, 191–207. https://doi.org/10.1016/j.mathsocsci.2016.09.001
- OpenAI. (2024). *How ChatGPT and our language models are developed*. https:// help.openai.com/en/articles/7842364-how-chatgpt-and-our-language-modelsare-developed
- Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C., Mishkin, P., Zhang, C., Agarwal, S., Slama, K., Ray, A., Schulman, J., Hilton, J., Kelton, F., Miller, L., Simens, M., Askell, A., Welinder, P., Christiano, P., Leike, J., & Lowe, R. (2022). *Training language models to follow instructions with human feedback*. https://proceedings.neurips.cc/paper_files/paper/2022/file/b1efde53be364a73914f58805a001731-Paper-Conference.pdf

- Phelps, S., & Russell, Y. I. (2023). Investigating emergent goal-like behaviour in large language models using experimental economics. https://doi.org/10.48550/ arXiv.2305.07970
- Sclar, M., Choi, Y., Tsvetkov, Y., & Suhr, A. (2023). Quantifying language models' sensitivity to spurious features in prompt design or: How i learned to start worrying about prompt formatting. https://doi.org/10.48550/arXiv.2310.11324
- Strachan, J. W. A, Albergo, D., Borghini, G., Pansardi, O., Scaliti, E., Gupta, S., Saxena, K., Rufo, A., Panzeri, S., Manzi, G., Graziano, M. S. A., & Becchio, C. (2024). Testing theory of mind in large language models and humans. *Nature Human Behaviour*, 8(7), 1285–1295. https://doi.org/10.1038/s41562-024-01882-z
- Trality. (2024). *Crypto trading bots: The ultimate beginner's guide*. Retrieved January 23, 2024 from https://medium.com/trality/crypto-trading-bots-f46405a7be11
- Tversky, A., & Kahneman, D. (1981). The framing of decisions and the psychology of choice. *Science*, *211*(4481), 453–458. https://www.jstor.org/stable/1685855
- Webb, T., Holyoak, K. J., & Lu, H. (2023). Emergent analogical reasoning in large language models. *Nature Human Behaviour*, 7(9), 1526–1541. https://doi. org/10.1038/s41562-023-01659-w

Aims and Scope

The aim of **Economics and Business Review** is to provide a platform for academicians from all over the world to share, discuss and integrate their research in the fields of economics and finance, including both behavioural economics and finance, with a key interest in topics that are relevant for emerging market economies. The journal welcomes submissions of articles dealing with micro, mezzo and macro issues that are well founded in modern theories or based on empirical studies and which are valuable for an international readership.

Your paper your way policy

The authors are initially expected to adjust their manuscripts to meet the basic requirements presented in the submission checklist below. Once the text has been accepted for publication, authors must adhere to all guidelines available on our website: https://journals.ue.poznan.pl/ebr

Basic requirements

- The submission has not been previously published nor is it under consideration for publication elsewhere (or an explanation has been provided in Comments to the Editor).
- The submitted manuscript must be anonymous. A separate title page must also be submitted, specifying each author's affiliation, email address, and ORCID iD. Acknowledgements and references to research grants should be included on the title page.
- The manuscript should be prepared in OpenOffice, Microsoft Word, or RTF document file format.
- The length of the manuscript should not exceed 8,000 words (including figures and tables, references, and footnotes).
- The manuscript includes an abstract of 100 to 150 words and is divided into clearly distinctive sections, including Introduction and Conclusions. The Introduction should state the aim of the study, research methods, main results, and particularly the study's contribution to international literature. The final paragraph should outline the article's content.
- All tables and figures should be numbered and presented consecutively according to their order in the text. Tables and figures should be as self-explanatory as possible, so readers do not need to refer to the main text to understand the information presented. The sources of all data used in tables and figures must be specified.
- The authors should use a consistent referencing style throughout the text.

The submission must be made via the submission system: https://journals.ue.poznan.pl/ebr/submission

More information and advice on the suitability and formats of manuscripts can be obtained from:

Economics and Business Review al. Niepodległości 10 61-875 Poznań Poland e-mail: secretary@ebr.edu.pl https://journals.ue.poznan.pl/ebr

Subscription

Economics and Business Review (EBR) is published quarterly and is the successor to the Poznań University of Economics Review. The EBR is published by the Poznań University of Economics and Business Press.

Economics and Business Review is indexed and distributed in Scopus, Claritave Analytics, DOAJ, ERIH plus, ProQuest, EBSCO, CEJSH, BazEcon, Index Copernicus and De Gruyter Open (Sciendo).