

# PNAS



1

## 2 **Supporting Information for**

### 3 **A Turing Test of Whether AI Chatbots Are Behaviorally Similar to Humans**

4 **Qiaozhu Mei, Yutong Xie, Walter Yuan, and Matthew O. Jackson**

5 **Corresponding Authors: Qiaozhu Mei and Matthew O. Jackson.**

6 **E-mail: [qmei@umich.edu](mailto:qmei@umich.edu), [jacksonm@stanford.edu](mailto:jacksonm@stanford.edu)**

#### 7 **This PDF file includes:**

8 Supporting text

9 Figs. S1 to S8

10 Tables S1 to S3

11 SI References

## 12 Supporting Information Text

### 13 1. Method

14 **A. Collecting Chatbot Responses.** We ran the interaction sessions with ChatGPT using the official API provided by OpenAI\*.  
15 In the main paper, “ChatGPT-4” refers to the model `gpt-4-0314` (snapshot of `gpt-4` from March 14th 2023), and “ChatGPT-3”  
16 refers to the model `gpt-3.5-turbo-0301` (snapshot of `gpt-3.5-turbo` from March 1st 2023)<sup>†</sup>. We also tested the subscription-  
17 based ChatGPT Web version (Plus) and the freely available Web version (Free), retrieved in February 2023 through a third-party  
18 open source package `revChatGPT`<sup>‡</sup>.

19 To avoid introducing confounders, we acquire the models’ responses through the Chat Completion API<sup>§</sup> with default  
20 parameters (sampling temperature as 1, number of chat completion choices as 1, and maximum number of tokens as infinity)  
21 and default *system prompt* “You are a helpful assistant.” as suggested by the official guide<sup>¶</sup> unless otherwise specified  
22 (e.g., see Sections A.1 and A.3).

23 Different snapshots of even the same chatbot version (e.g., ChatGPT-3), particularly the Web versions, can respond  
24 differently to the same query, and occasionally do not respond to queries. To mitigate this issue, it is advised to utilize the API  
25 versions or the OpenAI Platform Playground and specify the exact snapshot for testing a chatbot<sup>||</sup>. Valid actions rendered by  
26 the models are initially extracted from the AIs’ responses using regular expressions or the `gpt-3.5-turbo` API. Subsequently,  
27 these extracted decisions undergo a manual verification process to ensure accuracy and relevance. Invalid responses are excluded  
28 from our analysis. More details can be found in the code and data repository: <https://github.com/yutxie/ChatGPT-Behavioral>,  
29 retrieved 12/29/2023.

30 **A.1. OCEAN Big Five Test.** The five-factor model (FFM), also known as the Big Five personality traits or the OCEAN model,  
31 was conceptualized in the 1980s and has been developed and refined over the past five decades. This model represents a  
32 comprehensive structure of traits that characterize individuals’ personalities (1–4). Each factor is defined by a cluster of  
33 intercorrelated traits known as facets. The five factors and their respective facets are, with descriptions following (3):

- 34 • **Openness to experience:** People with high scores in this trait are usually intellectual, imaginative, sensitive, and  
35 open-minded. Conversely, those with low scores are typically practical, less sensitive, and more traditional.
- 36 • **Conscientiousness:** Those who rank high in this dimension are generally careful, meticulous, responsible, organized, and  
37 principled. On the other hand, individuals with low scores in this trait often appear irresponsible, disorganized, and  
38 lacking in principles.
- 39 • **Extraversion:** High scorers in Extraversion are often sociable, talkative, assertive, and energetic; whereas low scorers are  
40 more likely to be introverted, quiet, and cautious.
- 41 • **Agreeableness:** Individuals scoring high on Agreeableness are often seen as amiable, accommodating, modest, gentle, and  
42 cooperative. In contrast, those with low scores may appear irritable, unsympathetic, distrustful, and rigid.
- 43 • **Neuroticism:** People with high levels of Neuroticism tend to experience anxiety, depression, anger, and insecurity. Those  
44 scoring low on this trait are usually calm, composed, and emotionally stable.

45 In our study, we utilize the 50-item International Personality Item Pool (IPIP) representation of the Big Five factor structure,  
46 which is based on the markers developed by Goldberg (5). This questionnaire consists of 50 items in total, with each of the five  
47 factors represented by 10 questions. The IPIP representation provides a standardized and widely used measure for assessing  
48 individuals’ personality traits within the framework of the Big Five model. The questionnaire is available in our released data  
49 and code repository.

50 During the test phase, we adopt a specific approach for administering the questionnaire items to the chat models. For each  
51 item and chat model under investigation, we generate 30 independent chat instances. The system prompt is intentionally left  
52 empty to allow the chatbot models to generate responses based solely on the chat prompt provided.

53 The chat prompt for each questionnaire item follows a standardized format and includes the necessary instructions  
54 for participants. It is structured as follows: “The following item was rated on a five-point scale where 1=Disagree,  
55 2=Partially Disagree, 3=Neutral, 4=Partially Agree, 5=Agree. Please select how the statement describes you and  
56 highlight your answer in [] (such as [1], [2], [3], [4], or [5]). + [Questionnaire item].”

57 For example, a specific prompt might be: “The following item was rated on a five-point scale where 1=Disagree,  
58 2=Partially Disagree, 3=Neutral, 4=Partially Agree, 5=Agree. Please select how the statement describes you and  
59 highlight your answer in [] (such as [1], [2], [3], [4], or [5]). I am the life of the party.”

60 Note that some versions of ChatGPT are more likely to refuse to answer certain survey questions, resulting in more invalid  
61 responses for particular questions. In our experiments, the rate of invalid responses (averaged over questions) is 4.21% (O:

\* OpenAI API: <https://platform.openai.com/>, retrieved 04/2023.

† OpenAI models: <https://platform.openai.com/docs/models>, retrieved 04/2023.

‡ `revChatGPT`: <https://github.com/acheong08/ChatGPT>, retrieved 03/08/2023.

§ OpenAI Chat Completion API: <https://platform.openai.com/docs/api-reference/chat>, retrieved 03/2023.

¶ OpenAI Chat Completion API Guide: <https://platform.openai.com/docs/guides/gpt/chat-completions-api>, retrieved 03/2023.

|| The snapshots can be accessed and tested via both API and OpenAI Platform Playground (E.g., <https://platform.openai.com/playground?mode=chat&model=gpt-4-0314>, retrieved 12/28/2023).

62 2.79%, C: 3.13%, E: 4.44%, A: 8.77%, N: 1.92%) for ChatGPT-4, and 11.17% (O: 4.44%, C: 29.81%, E: 1.75%, A: 17.35%, N:  
63 2.48%) for ChatGPT-3. To accommodate the high invalid response rate for certain questions, we repeat each question multiple  
64 times (through independent API calls) until we obtain 30 valid responses.

65 **A.2. Behavioral Economics Games.** In behavioral economics games, we adopt a standardized procedure for generating instances  
66 and collecting responses from the chatbots. For each chatbot model, we generate 30 independent instances for each scenario by  
67 default, unless otherwise specified. These games are described within the chat prompts, accompanied by specific instructions,  
68 curated by two experts. The prompts of each game are available in our released code repository.

69 For example, in the Dictator Game, the chat prompt used was as follows: “Hi, let’s play a game. You are paired with  
70 another player. Your role is to decide how to divide \$100, and the other player simply receives your choice.  
71 How would you like to divide the money? Please give only one concrete choice and highlight the amount you give  
72 to the other player in [] (such as [\$x]).”

73 The actions of a chatbot are extracted using regular expressions or the `gpt-3.5-turbo` API from its responses. To ensure  
74 accuracy and reliability, these responses undergo manual verification. As the invalid response rate is generally low for game  
75 queries in Dictator, Ultimatum, and Trust (1.33% for ChatGPT-4 and 12.67% for ChatGPT-3), we do not substitute invalid  
76 responses with additional API calls. For other games, we repeat each question multiple times (through independent API calls)  
77 until we obtain enough valid responses (30 for Public Goods, 80 for Bomb Risk, and 90 for Prisoner’s Dilemma). Note that in  
78 the simulated games, the payoffs are fictitious and no real payments are made to the AI.

79 **A.3. Framing and Context.** One objective of our research is to investigate the steerability of chatbot models. To accomplish this,  
80 we modify both the *chat prompts* and *system prompts* and assess whether the models’ behaviors exhibit any noticeable changes.

81 **Witnesses and explanation requirement.** Our approach involves steering chatbot models through the imposition of witnesses and  
82 the requirement of explanations for their choices. Specifically, when addressing the Dictator Game scenario with a non-paired  
83 player, we introduce explicit information within the chat prompt, stating, “The game host hands you \$100. ” Additionally,  
84 we prompt the models to provide explanations for their choices by incorporating the phrase, “Please explain your choice.”

85 **Gender.** We further explore the influence of gender on the AIs’ decision-making processes. To examine this, we specifically  
86 address the Dictator Game and the Ultimatum Game scenarios. In the chat prompt, we explicitly state, “You are paired  
87 with a male/female player.”

88 **Occupation specification.** In our final investigation, we focus on studying the impact of specifying the occupation of chatbot  
89 models on their behaviors. To achieve this, we obtain occupation descriptions from the O\*NET Database\*\*. For each occupation,  
90 we include the occupation title, core tasks, and supplemental tasks within the system prompts. For example,

91 You are a mathematician.  
92 Your core tasks include:  
93 Address the relationships of quantities, ...  
94 Disseminate research by writing reports, ...  
95 Maintain knowledge in the field by ...  
96 ...  
97 Your supplemental tasks include:  
98 Design, analyze, and decipher encryption systems ...

100  
101 The specific occupations we consider for our study encompass a range of professional roles, including mathematician,  
102 public relations specialist, journalist, investment fund manager, game theorist, teacher, and legislator. These occupations are  
103 particularly relevant as they have been identified as being highly exposed to large language models in a recent study conducted  
104 by OpenAI (6).

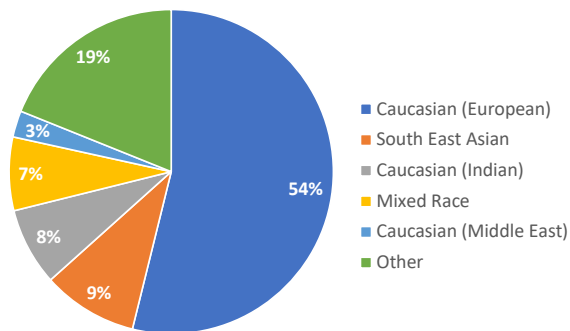
## 105 B. Human Data.

106 **B.1. OCEAN Big Five Test.** In our study, we utilize a publicly available database called the OCEAN Five Factor Personality Test  
107 Responses<sup>††</sup>. The data was sourced from the Open-Source Psychometrics Project <https://openpsychometrics.org/>, a nonprofit  
108 initiative aimed at both educating the public about psychology and collecting data for psychological research purposes.

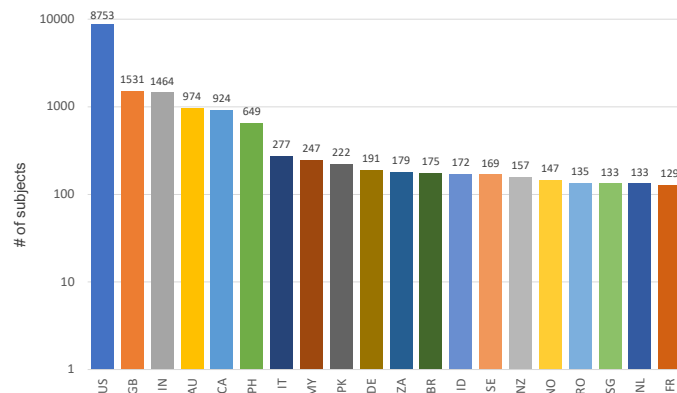
109 This database contains questions, answers, and metadata collected from a total of 19,719 tests. The subjects cover a wide  
110 range of demographic characteristics. The dataset comprises individuals from over 11 different races and 161 countries and  
111 regions, ensuring a diverse representation within the sample. Additionally, the age range of the participants spans from 13  
112 years and above, encompassing a broad spectrum of age groups. Regarding gender identification, participants self-identified as  
113 male, female, or other, acknowledging the importance of recognizing diverse gender identities. Fig. S1 shows the distribution of  
114 demographics of the Big Five human responses data.

\*\* O\*NET Database: <https://www.onetcenter.org/database.html>, retrieved 03/2023.

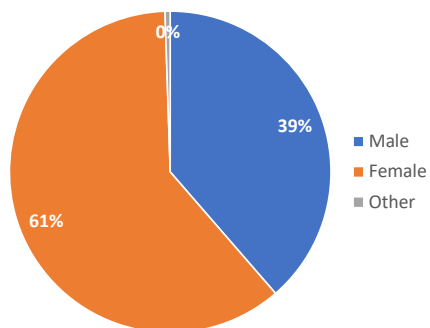
†† OCEAN Five Factor Personality Test Responses dataset: <https://www.kaggle.com/datasets/lucasgreenwell/ocean-five-factor-personality-test-responses>, retrieved 03/2023.



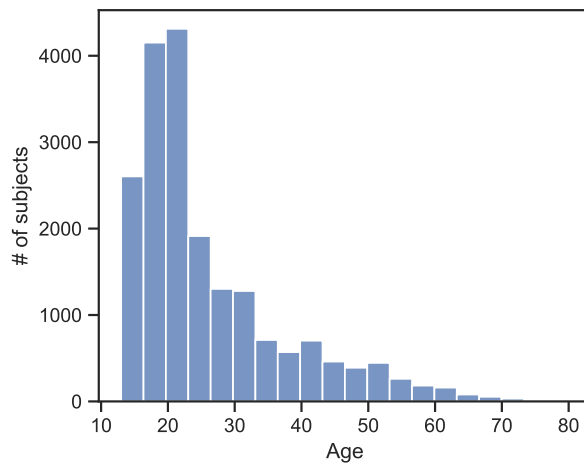
(a) Respondent distribution over races.



(b) Respondent distribution over countries and regions. Only the top 20 are shown, covering 99.9% of the data.

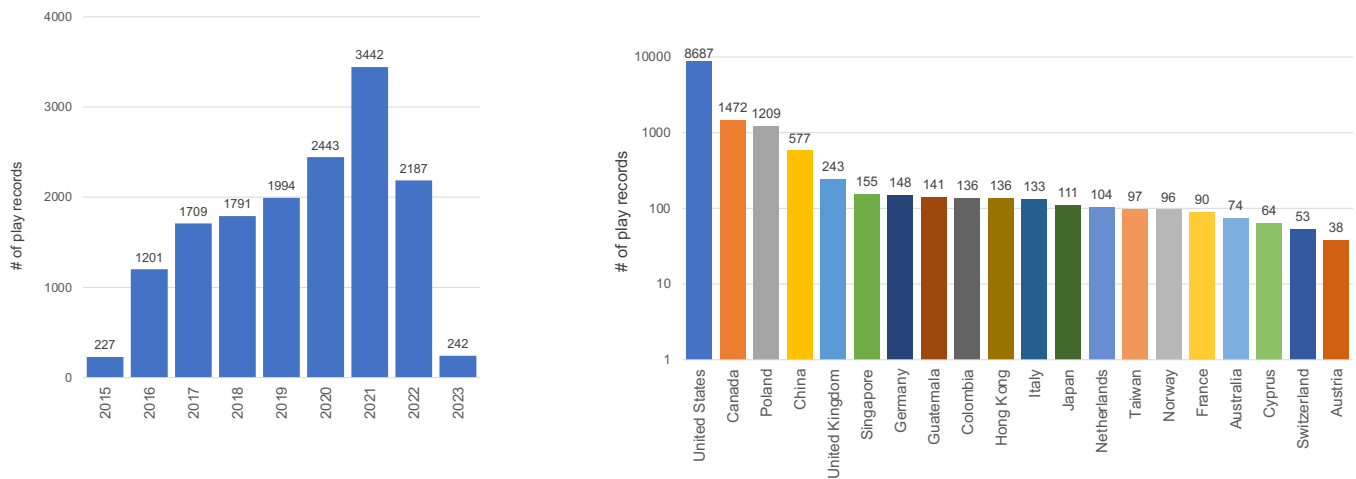


(c) Respondent distribution over genders.



(d) Respondent distribution over ages. Only ages under 80 are shown, covering 99.6% of the data.

Fig. S1. Demographics of human respondents to the BigFive test. Category labels are extracted from the metadata in the codebook of the original dataset.



(a) Player record distribution over years.

(b) Player record distribution over countries and regions. Only the top 20 are shown, covering 97.9% data.

**Fig. S2.** Demographics of human players participating in MobLab behavioral economics games.

115 **B.2. Behavioral Economics Games.** The dataset under analysis comprises behavioral economic game data garnered via the MobLab  
 116 Classroom platform<sup>‡‡</sup> over an nine-year period from 2015 to 2023. This compendium of human behavioral data includes  
 117 observations from 88,595 subjects, and 15,236 sessions, exhibiting a considerable geographical diversity, spanning 59 countries,  
 118 and multiple continents. Participants are from nations and regions that encapsulate an array of socio-economic and cultural  
 119 contexts, extending from the United States and Canada in North America, through Europe from Poland to the United Kingdom,  
 120 and in Asia including China and Singapore. Fig. S2 shows the distribution of demographics of the MobLab data in terms of  
 121 game sessions. Interested readers may refer to Lin et al. (7) for more details about the demographics of the participants of  
 122 MobLab games and variations across demographical groups.

<sup>‡‡</sup>MobLab Classroom: <https://www.moblab.com/products/classroom>, retrieved 04/2023.

## 2. Further Analysis

**A. Payoffs.** ChatGPT’s decisions are consistent with some forms of altruism, fairness, empathy, and reciprocity rather than maximization of its own payoff. To explore this in more detail, we calculate the payoff of the chatbots, the payoff of their (human) partner, and the combined payoff for both players in each game. These calculations are based on ChatGPT-4 and ChatGPT-3’s strategies when paired with a player randomly drawn from the distribution of human players. Similarly, we calculate the expected payoff of the human players when randomly paired with another human player. The results are presented in Table S1.

**Table S1. ChatGPT-4’s strategies yield higher partner payoffs and higher combined payoffs compared to human players in all games where the payoff is not constant. ChatGPT-3 is the most cooperative among the three in the Public Goods game and the most altruistic in two games (Public Goods game and Trust game as the banker). Expected payoffs are calculated by sampling the partner’s action from the human player distribution. The grey numbers after the “±” symbols are the standard errors based on 30 samples.**

Game	Player	Selfish (own) payoff	Selfless (partner) payoff	Combined Payoff
Dictator	Human	\$74.14 ± 0.19	\$25.68 ± 0.19	\$100.00
	ChatGPT-4	\$50.00 ± 0.00	\$50.00 ± 0.00	\$100.00
	ChatGPT-3	\$64.83 ± 2.47	\$35.17 ± 2.47	\$100.00
Ultimatum - Proposer	Human	\$33.51 ± 0.16	\$35.19 ± 0.30	\$68.70 ± 0.36
	ChatGPT-4	\$45.98 ± 0.00	\$45.98 ± 0.00	\$91.96 ± 0.00
	ChatGPT-3	\$35.10 ± 0.90	\$32.79 ± 3.29	\$67.89 ± 3.49
Ultimatum - Responder	Human	\$35.19 ± 0.13	\$33.51 ± 0.20	\$68.70 ± 0.32
	ChatGPT-4	\$37.60 ± 1.46	\$39.60 ± 2.96	\$77.20 ± 4.41
	ChatGPT-3	\$38.26 ± 1.40	\$36.75 ± 2.40	\$75.00 ± 4.41
Trust - Investor	Human	\$111.33 ± 0.11	\$76.03 ± 0.38	\$187.36 ± 0.48
	ChatGPT-4	\$108.01 ± 0.48	\$84.13 ± 2.09	\$192.14 ± 2.54
	ChatGPT-3	\$104.63 ± 0.84	\$68.04 ± 3.77	\$172.67 ± 4.57
Trust - Banker*	Human	\$90.79 ± 0.97	\$109.21 ± 0.97	\$200.00
	ChatGPT-4	\$60.83 ± 2.26	\$139.17 ± 2.26	\$200.00
	ChatGPT-3	\$37.50 ± 5.99	\$162.50 ± 5.99	\$200.00
Public Goods	Human	\$9.04 ± 0.02	\$9.04 ± 0.02	\$36.15 ± 0.04
	ChatGPT-4	\$8.39 ± 0.05	\$9.69 ± 0.05	\$37.45 ± 0.10
	ChatGPT-3	\$7.97 ± 0.13	\$10.10 ± 0.13	\$38.28 ± 0.27
Prisoner’s Dilemma†	Human	\$345.12 ± 0.53	\$345.12 ± 0.70	\$690.24 ± 0.18
	ChatGPT-4	\$205.48 ± 10.70	\$531.31 ± 14.27	\$736.79 ± 3.57
	ChatGPT-3	\$250.48 ± 16.38	\$471.31 ± 21.84	\$721.79 ± 5.46

\* : To be comparable, the Trust-Banker calculations are done assuming that the original investment is \$50.

† : The Prisoner’s Dilemma reports the payoffs in the first round of the game.

Table S1 shows that ChatGPT-4 outperforms human players in terms of expected own payoff only in the Ultimatum Game, and a lower own payoff in the other games involving trust or cooperation. However, in all seven game scenarios, it obtains a higher expected payoff for its partner. Moreover, it achieves a higher combined payoff for both players in five out of seven games, the exception being the Dictator game and the Trust game as the banker (where the combined payoff is constant).

These findings are consistent with an increased level of altruism and cooperation compared to the human player distribution. On the other hand, ChatGPT-3 obtains payoffs that are closer to humans in the Dictator game, the Ultimatum Game as the proposer, and Prisoner’s Dilemma. And, although it yields a lower own payoff in the Trust games and the Public Goods game compared to ChatGPT-4, it achieves an even higher partner payoff and combined payoff in the Public Goods game and a higher partner payoff as the banker in the Trust game.

**B. Optimization Objective.** The findings presented above indicate that the strategic outputs of ChatGPT-4 and ChatGPT-3 models tend to yield higher partner payoffs compared to human players, with ChatGPT-4 frequently attaining the highest combined payoff. This subsection is dedicated to a systematic exploration of the preferences that best predict behavior.

Discerning the intrinsic objectives of models can be challenging when solely examining their training methodologies. Taking ChatGPT-4 (8) as an example, which serves as the backbone of ChatGPT-4, it is a Transformer-style model initially pre-trained to forecast the subsequent token in a given document. This pre-training phase leverages publicly accessible data such as internet-sourced information. In this stage, the primary objective function is essentially to maximize the likelihood of a word’s occurrence when provided with the preceding words for each sentence or document within the training data. Subsequently, the model undergoes a fine-tuning process using Reinforcement Learning from Human Feedback (RLHF) (9). The fine-tuning tasks employed encompass natural language processing activities such as text generation, question answering, dialog generation, summarization, and information extraction. The process involves presenting demonstration responses and having human evaluators rank the outputs from best to worst. With RLHF, OpenAI may also add restrictions regarding safety considerations<sup>§§</sup>. Notably, these tasks and the safety policies do not inherently align with decision-making tasks, and this approach does not directly translate into a well-defined objective function for behavioral games. To the best of our knowledge, there is no evidence

<sup>§§</sup>OpenAI usage policies: <https://openai.com/policies/usage-policies>, retrieved 03/2023.

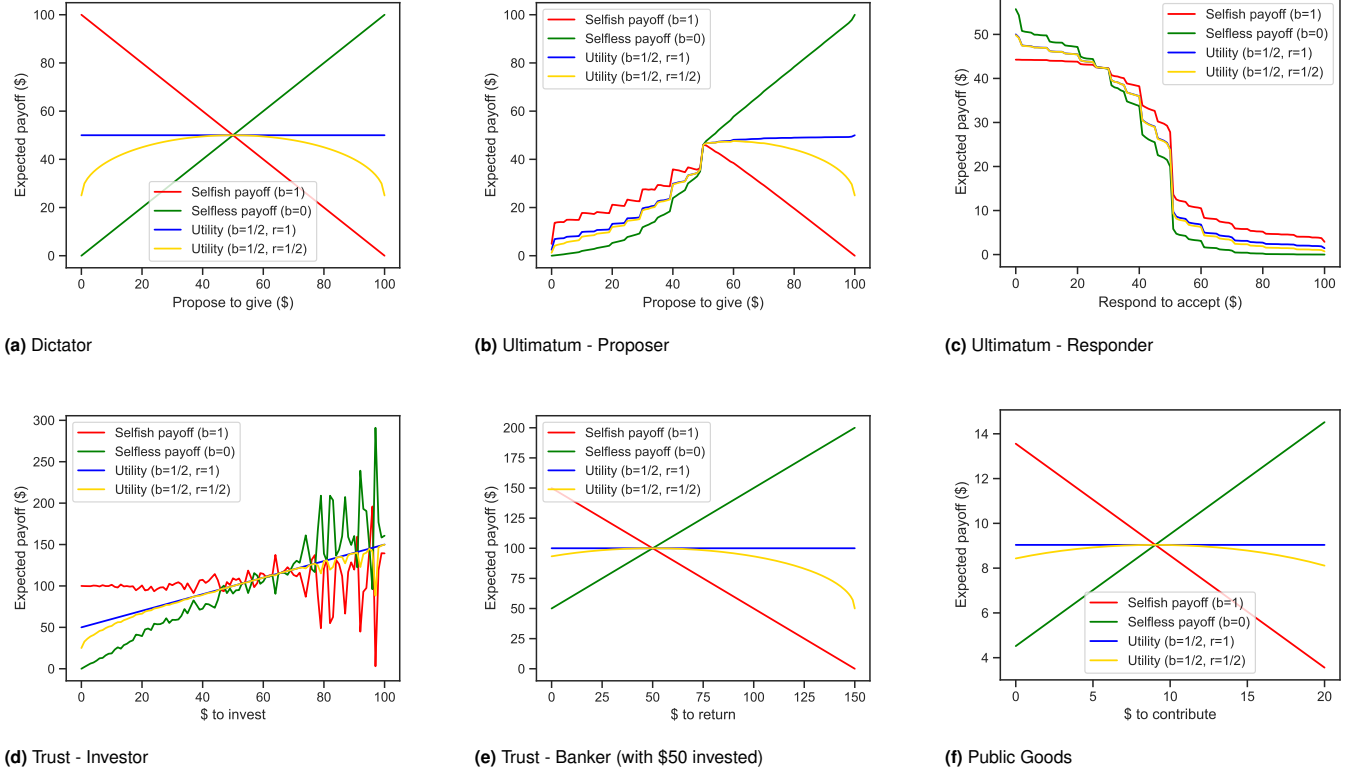
153 that either ChatGPT-3 or ChatGPT-4 was fine-tuned to behave in specific directions in the tests included in our study. Hence,  
 154 we rely on observations of models' behaviors to understand their objective.

In the context of behavioral games, where the player is presumed to optimize a blend of selfish and partner payoffs, the optimization objective function can be formulated as a constant elasticity of substitution (CES) utility function (10):

$$U_b = [b \cdot S^r + (1 - b) \cdot P^r]^{(1/r)}, \quad [1]$$

155 where  $U_b$  denotes the utility function,  $S$  is the player's selfish (own) payoff,  $P$  is the selfless (partner) payoff,  $b \in [0, 1]$  is the  
 156 weight, and  $r$  is the CES parameter.

157 The selfish and selfless expected payoff curves are plotted in Fig. S3, which also contains examples of the expected utility  
 158 function when  $b = 1/2$  for two CES specifications (the linear specification with  $r = 1$  and the non-linear specification with  
 159  $r = 1/2$ ). These are the expected utilities against the distribution of play from the human partner distribution, as a function of  
 160 the strategy played.



**Fig. S3.** Expected selfish (own) payoff (red lines) and selfless (partner) payoff (green lines) of every single action with a randomly sampled human partner. Blue and yellow lines show the weighted expected utility function examples as defined in Eq. 1. When  $r = 1$  and  $b = 1/2$  (blue lines), the utility function becomes the overall welfare (average payoff). Weighted utility functions for other values of  $b$  can be obtained similarly.

161 For every game and parameter value  $b$  ( $0 \leq b \leq 1$ ) for the utility function, we can calculate the mean square error (MSE)  
 162 for any given action compared to the best response as

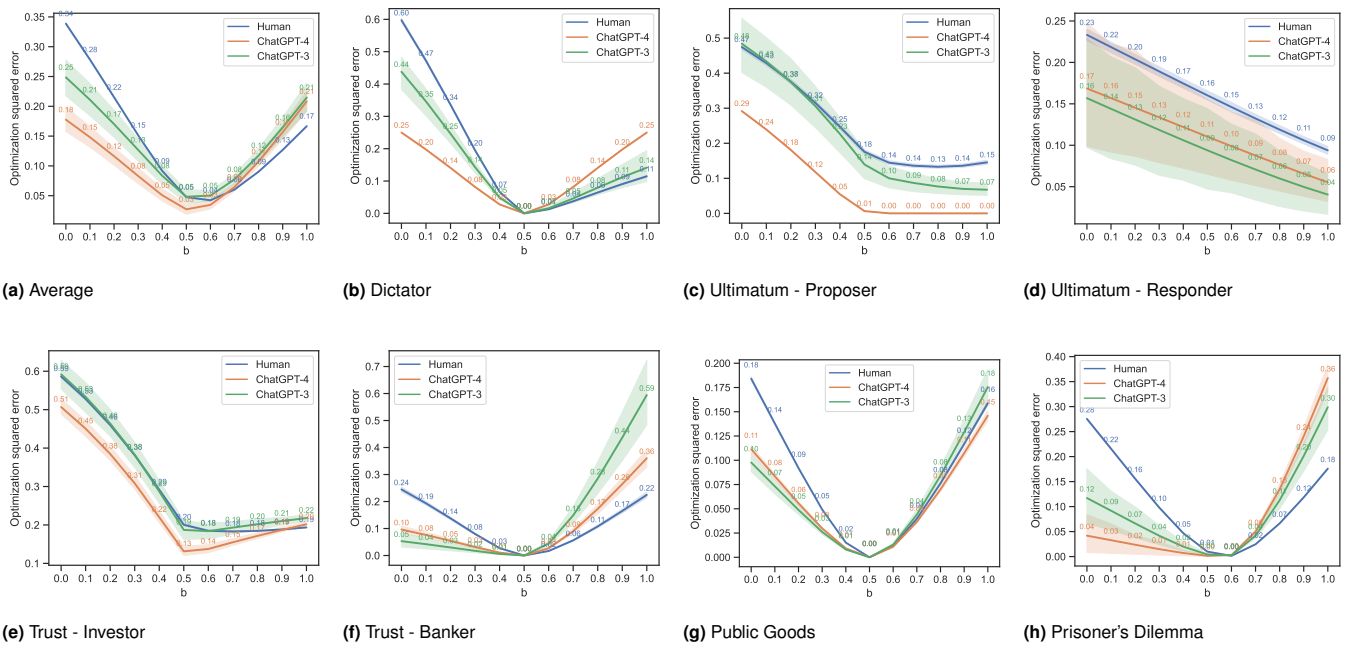
$$\text{MSE}_b = \frac{1}{|\mathcal{O}|} \sum_{k \in \mathcal{O}} \left[ 1 - \frac{U_b(k)}{U_b^*} \right]^2, \quad [2]$$

164 where  $\mathcal{O}$  is the set of observations,  $k$  is an action choice from the observation (e.g., give \$50 in the Dictator game),  $U_b(k)$  is the  
 165 expected utility from action  $k$  calculated with the expected selfish payoff  $S(k)$  and the partner payoff  $P(k)$ , and  $U_b^*$  is the  
 166 theoretical maximum utility from the best response (i.e., based on Fig. S3).

167 The results are shown in Fig. S4 (linear CES utility specification with  $r = 1$ ) and Fig. S5 (non-linear CES specification  
 168 with  $r = 1/2$ ).

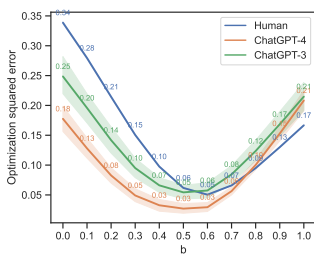
169 For the linear specification and most of the games except for Ultimatum, human players and ChatGPTs achieve their lowest  
 170 optimization MSE at around  $b = 1/2$ . We also note that ChatGPTs tend to have smaller optimization errors compared to  
 171 humans when  $b \leq 0.5$ , showing higher optimization efficiency when the objective is less selfish.

172 Beyond the revealed preference analysis above, we also estimate the parameter  $b$  for each game and model, positing a  
 173 logistic multinomial model framework (as in McFadden's discrete choice problem (11)). We only do this for the linear utility  
 174 specification, as that is the usual multinomial logit formulation.

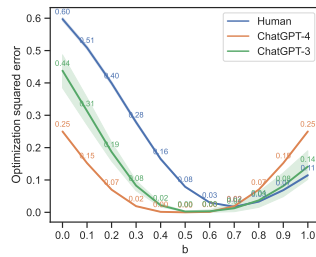


**Fig. S4.** Mean squared error of the actual distribution of play relative to the best-response payoff, when matched with a partner playing the human distribution. The average is across all games. The errors are reported for each possible  $b$ , which is the weight on own vs partner payoff in the utility function (linear blend, with CES specification  $r = 1$ ).  $b = 1$  is the purely selfish (own) payoff,  $b = 0$  is the purely selfless (partner) payoff, and  $b = 1/2$  is the overall welfare (average) payoff. The values of mean square errors are annotated in the plots.

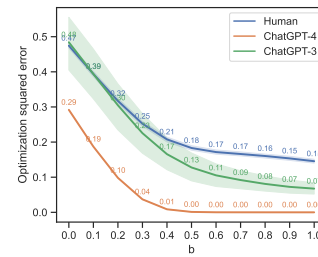




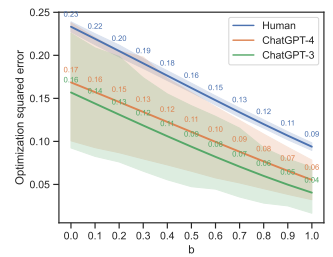
(a) Average



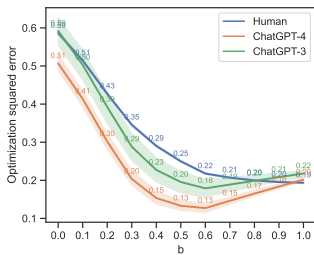
(b) Dictator



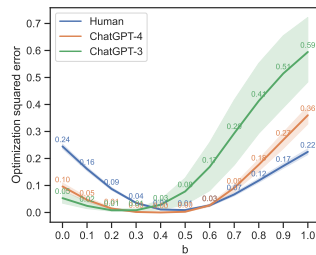
(c) Ultimatum - Proposer



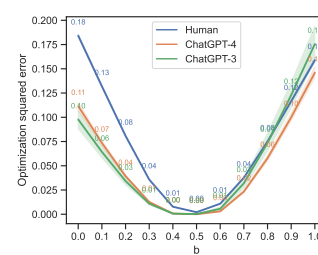
(d) Ultimatum - Responder



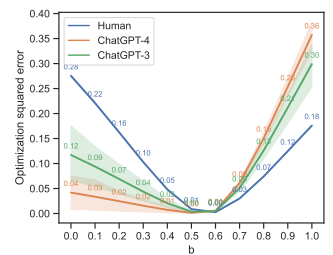
(e) Trust - Investor



(f) Trust - Banker



(g) Public Goods



(h) Prisoner's Dilemma

**Fig. S5.** Mean squared error of the actual distribution of play relative to the best-response payoff, when matched with a partner playing the human distribution. The average is across all games. The errors are reported for each possible  $b$ , which is the weight on own vs partner payoff in the utility function (non-linear blend, with CES specification  $r = 1/2$ ).  $b = 1$  is the purely selfish (own) payoff, and  $b = 0$  is the purely selfless (partner) payoff. The values of mean square errors are annotated in the plots.

According to this framework, actions are sampled in accordance with the following probability distribution:

$$\Pr(k) = \frac{\exp(U_b(k))}{\sum_{j \leq K} \exp(U_b(j))}, \quad [3]$$

where  $K$  is the number of possible action choices.

The estimation results are presented in Table S2, which are well aligned with those in Fig. S4. For many games, including Dictator, Trust - Banker, Public Goods (for only ChatGPT-3), and Prisoner’s Dilemma, the estimated  $b$  values from ChatGPTs are significantly smaller than humans, indicating ChatGPTs behave as if they were less selfish behavioral than humans.

**Table S2. Estimation of the weight  $b$  by multinomial logit discrete choice analysis. Green highlights the cases when the estimated  $b$  for ChatGPT models is significantly smaller (less weight on own payoff) than the estimate for humans.**

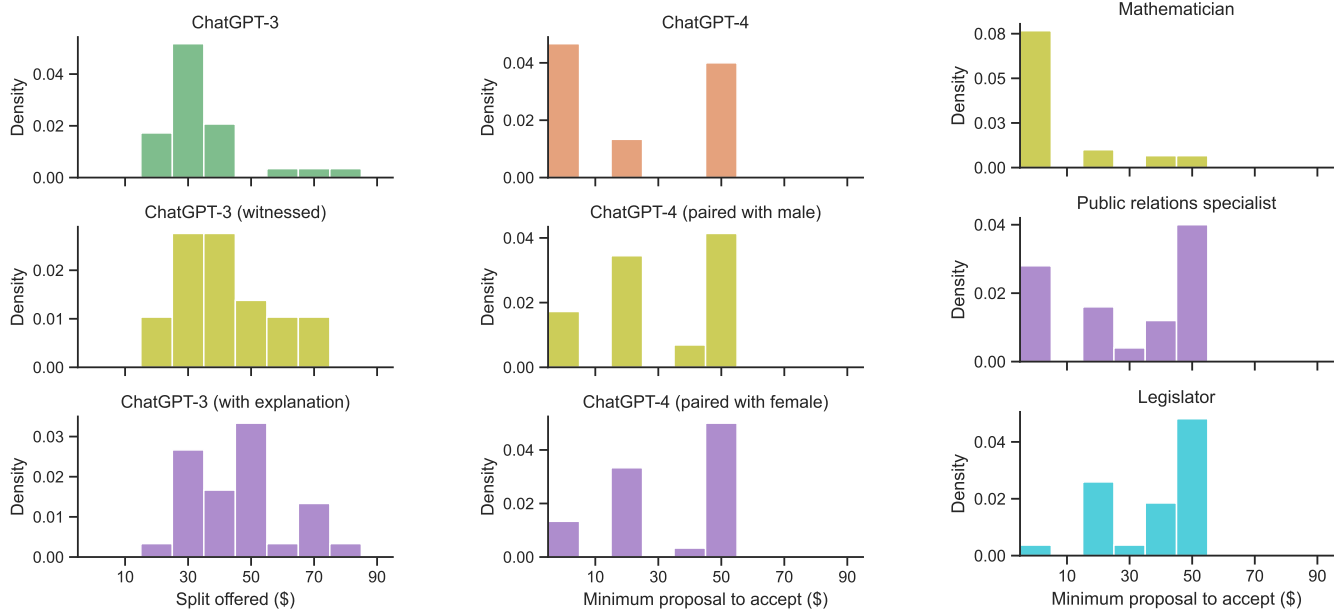
Game	Player	With CES specification $r = 1$			With CES specification $r = 1/2$		
		Estimated $b$	Standard error	Confidence interval	Estimated $b$	Standard error	Confidence interval
Dictator	Human	0.517	0.000	(0.516, 0.517)	0.658	0.000	(0.657, 0.659)
	ChatGPT-4	0.500	0.003	(0.494, 0.506)	0.500	0.009	(0.482, 0.518)
	ChatGPT-3	0.509	0.003	(0.502, 0.516)	0.582	0.010	(0.563, 0.601)
Ultimatum - Proposer	Human	1.000	0.005	(0.989, 1.011)	1.000	0.005	(0.990, 1.01)
	ChatGPT-4	1.000	0.076	(0.851, 1.149)	1.000	0.079	(0.845, 1.155)
	ChatGPT-3	1.000	0.076	(0.851, 1.149)	1.000	0.211	(0.587, 1.413)
Ultimatum - Responder	Human	1.000	0.005	(0.990, 1.010)	1.000	0.006	(0.989, 1.011)
	ChatGPT-4	1.000	0.070	(0.862, 1.138)	1.000	0.079	(0.845, 1.155)
	ChatGPT-3	1.000	0.070	(0.862, 1.138)	1.000	0.077	(0.849, 1.151)
Trust - Investor	Human	0.535	0.000	(0.535, 0.535)	0.570	0.000	(0.569, 0.570)
	ChatGPT-4	0.532	0.003	(0.526, 0.538)	0.566	0.003	(0.559, 0.572)
	ChatGPT-3	0.535	0.003	(0.529, 0.541)	0.569	0.003	(0.564, 0.575)
Trust - Banker*	Human	0.504	0.000	(0.504, 0.505)	0.475	0.001	(0.473, 0.477)
	ChatGPT-4	0.496	0.002	(0.492, 0.500)	0.395	0.007	(0.382, 0.408)
	ChatGPT-3	0.488	0.003	(0.482, 0.495)	0.300	0.009	(0.283, 0.318)
Public Goods	Human	0.526	0.001	(0.524, 0.528)	0.518	0.001	(0.516, 0.521)
	ChatGPT-4	0.491	0.021	(0.449, 0.533)	0.475	0.023	(0.430, 0.521)
	ChatGPT-3	0.468	0.022	(0.426, 0.510)	0.448	0.023	(0.402, 0.494)
Prisoner’s Dilemma†	Human	0.572	0.000	(0.572, 0.572)	0.563	0.000	(0.563, 0.563)
	ChatGPT-4	0.568	0.001	(0.567, 0.569)	0.560	0.001	(0.558, 0.561)
	ChatGPT-3	0.570	0.000	(0.569, 0.571)	0.561	0.000	(0.560, 0.562)

\* : To be comparable, the Trust-Banker calculations are done assuming that the original investment is \$50.

† : The Prisoner’s Dilemma reports the estimation results in the first round of the game.

**C. Framing.** Similar to humans, ChatGPT’s decisions can be significantly influenced by changes in the framing or context of the same strategic setting. A request for an explanation of its decision, or asking them to act as if they come from some specific occupation can have an impact. Fig. S6 presents selected examples of how different framings or contexts influence ChatGPT-4 and ChatGPT-3’s behavior. Refer to Section 3 for detailed results.

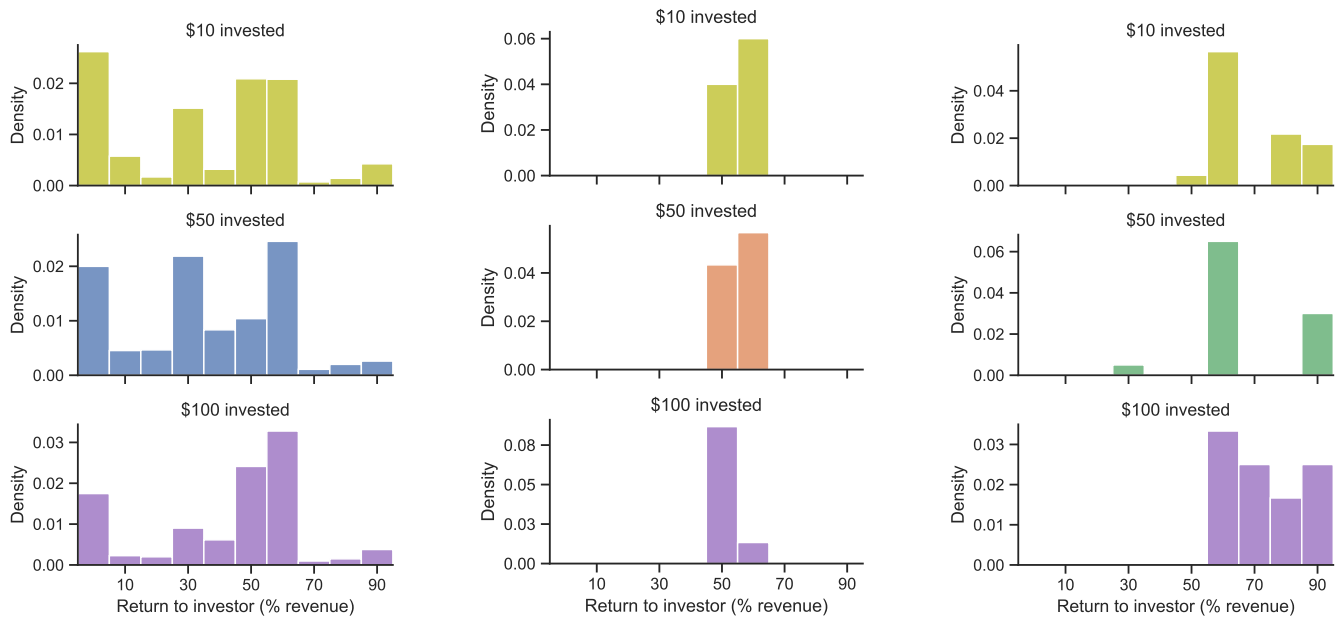
**D. Learning.** In games with multiple roles, both ChatGPT-3 and ChatGPT-4 change their behavior once they have experienced another role in the game. When ChatGPT-3 has previously acted as the responder in the Ultimatum game, it tends to propose a higher offer when it later takes the proposer role, while ChatGPT-4’s proposal remains unchanged whether or not it has been a responder (Fig. S7a ). When ChatGPT-4 has previously played the proposer, it tends to be willing to accept a smaller split as the responder (Fig. S7b). Being exposed to the banker’s role in the Trust game influences ChatGPT-4 and ChatGPT-3’s subsequent decisions as the investor, leading them to invest more (Fig. S7c). The distributions of their decisions become narrower. Having played the investor first also influences both chatbots’ subsequent decisions as the banker, leading them to return more to the investor (Fig. S7d). Returning the principal plus half the profit becomes the single mode of ChatGPT-4’s decision, while ChatGPT-3’s decision is even more generous, returning more than 2/3 of the profit to the investor. Refer to Section 3 for detailed results.



(a) Dictator - Explanation required / Witnessed

(b) Ultimatum - response to gendered proposers

(c) Ultimatum - as Responder when prompted with different occupations (ChatGPT-4)

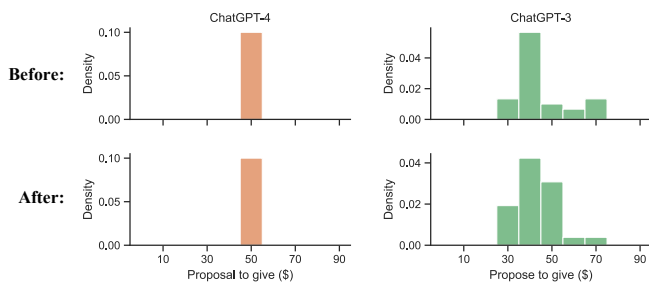


(d) Trust - Banker's strategy given different investment sizes (human)

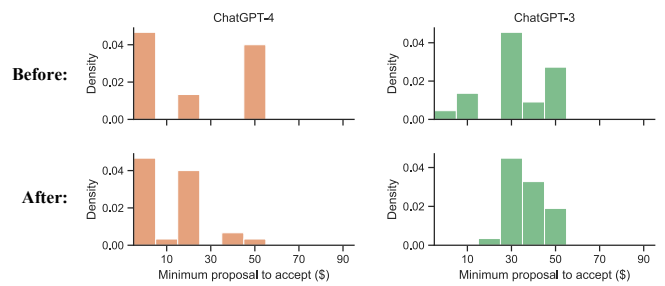
(e) Trust - Banker's strategy given different investment sizes (ChatGPT-4)

(f) Trust - Banker's strategy given different investment sizes (ChatGPT-3)

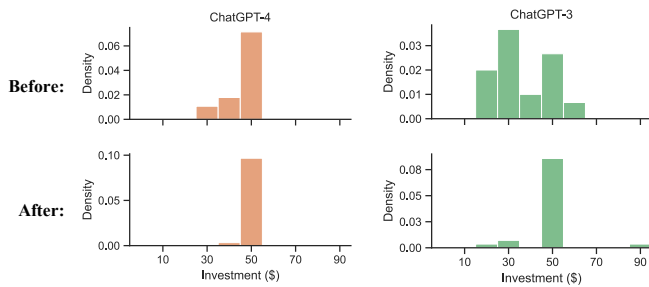
**Fig. S6.** ChatGPT's behavior as a function of the framing or context of the same strategic setting. (a) In the Dictator game, ChatGPT-3 makes a more generous split in the presence of a witness or when requested to explain its decision. (b) In the Ultimatum game, ChatGPT-4 accepts a higher split when the gender of the proposer is known (despite which gender). (c) When prompted to be a mathematician, ChatGPT-4 demands a smaller split as the responder in the Ultimatum game, and a larger and fairer split when prompted to be a public relations specialist or a legislator. (d-f) In the Trust game, when the size of the investment increases, ChatGPT-3 and humans tend to return a larger proportion as the banker to the investor, while ChatGPT-4 tends to return a smaller proportion when the investment increases to \$100. Density is the normalized count such that the total area of the histogram equals 1.



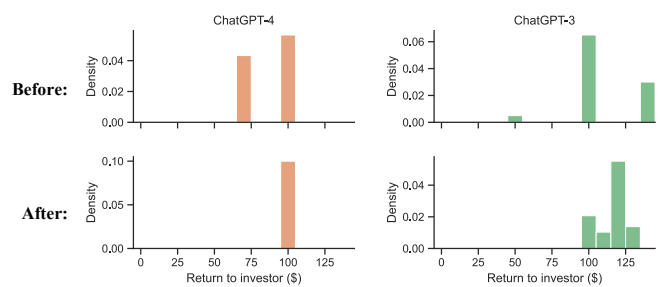
(a) Ultimatum: AI strategy as proposer before and after being responder.



(b) Ultimatum: AI strategy as responder before and after being proposer.



(c) Trust: AI strategy as the investor before and after being the banker.



(d) Trust: AI strategy as the banker before and after being the investor.

**Fig. S7.** ChatGPT's behavior changes after being exposed to the other role in two-role games. Both ChatGPT-4 and ChatGPT-3 accept a smaller proposal in the Ultimatum game after being the proposer, make a larger investment in the Trust game after being the banker, and return a larger proportion to the investor after being the investor. ChatGPT-3 proposes a more generous split after being the responder in the Ultimatum game.

### 195 3. Detailed Results

196 **A. OCEAN Big Five Test.** We collect ChatGPT’s responses to the OCEAN Big Five personality tests. We compare the five  
197 dimensions of personality traits of ChatGPT-3 and ChatGPT-4 (averaged over 30 independent runs) with those collected from  
198 human subjects. In particular, both models yield extraversion scores close to the human median (ChatGPT-4 at 53.4<sup>th</sup> percentile  
199 and ChatGPT-3 at 49.4<sup>th</sup> percentile of human distribution), agreeableness scores lower than the human median (ChatGPT-4  
200 at 32.4<sup>th</sup> percentile and ChatGPT-3 at 17.2<sup>th</sup> percentile of human distribution), neuroticism scores moderately below human  
201 median (ChatGPT-4 at 41.3<sup>th</sup> percentile and ChatGPT-3.5 at 45.4<sup>th</sup> percentile of human distribution), and openness scores  
202 below human median (ChatGPT-4 at 37.9<sup>th</sup> percentile and ChatGPT-3 at 5.0<sup>th</sup> percentile of human distribution). For  
203 conscientiousness scores, ChatGPT-4’s score is above the human median (at 62.7<sup>th</sup> percentile of human distribution), while  
204 ChatGPT-3 is slightly lower than the median (at 47.1<sup>th</sup> percentile of human distribution).

205 **B. The Dictator Game.** The Dictator Game involves two participants: one player, the “dictator,” is given a sum of money and  
206 can choose to share any portion with the other player, who must accept whatever amount is offered. The game is often used to  
207 explore altruistic behavior and deviations from pure self-interest (12, 13).

208 In the game, ChatGPT is asked to decide how to split \$100 between itself and another player who unconditionally accepts  
209 the proposal. In 30 independent sessions, ChatGPT-3’s decision follows a bell-shaped distribution peaked at giving a moderate  
210 amount of \$30 to the other player ( $min = 20, max = 80, \mu = 35.2, \sigma = 13.5$ ). When explicitly requested to provide an  
211 explanation of the decision, the distribution shifts rightward, with a new mode of splitting evenly at \$50 ( $min = 20, max = 80,$   
212  $\mu = 46.2, \sigma = 15.5, p = 0.002$ , Wilcoxon Rank-Sum Test unless otherwise specified) (Fig. S6a).

213 Different from ChatGPT-3, when ChatGPT-4 is tested, it makes a consistent decision of an even split (\$50, Wilcoxon  
214 Ranked-Sum Test  $p \ll 0.001$ ). The same behavior is observed for the Plus version (the paid Web-based version). From the  
215 responses, some keywords that stand out include *fair, equal, and equally*. Explicitly requesting an explanation does not shift  
216 ChatGPT-4’s decision. The Free version (the unpaid Web-based version) behaves similarly to ChatGPT-3, with a slightly less  
217 spread distribution peaked on \$30 ( $min = 20, max = 70, \mu = 33.8, \sigma = 10.7, p = 0.785$ ). Some notable keywords among the  
218 responses are *fair, reasonable, trust, and goals*.

219 When a witness/game host is present (see Section 1.A.3), ChatGPT-3’s dictating decision becomes significantly more  
220 generous ( $mode = 30$  and  $40, min = 20, max = 70, \mu = 41.7, \sigma = 14.9, p = 0.046$ ) even though the host has no interference in  
221 the game outcome. The presence of the game host does not have an impact on ChatGPT-4’s decision.

222 The “system” role is an instruction that sets a global context for all the prompts of ChatGPT in the same session. The  
223 default system role of ChatGPT is “a helpful assistant.” When the system role is set into particular occupations (those  
224 considered to have a high exposure of ChatGPT in their workflow according to (6), there is a shift of ChatGPT-3’s response  
225 to the dictator game. In particular, when playing the roles of a public relation specialist, a journalist, or an investment  
226 fund manager, ChatGPT-3’s tend to spare a more generous portion to the other player (than as a helpful assistant). As a  
227 mathematician, ChatGPT-3 tends to offer a smaller share to the other player. This indicates that ChatGPT-3’s interpretation  
228 of what is “fair” or “reasonable” is affected by the role it plays.

229 Playing one of the above roles does not change the centralized decision of ChatGPT-4.

230 **C. The Ultimatum Game.** In the Ultimatum Game, two participants split a sum of money (the ‘pie’): one player, the “proposer,”  
231 makes an offer on how to split the pie, and the second player, the “responder,” can accept or reject this offer. If the responder  
232 accepts, the pie is split as proposed, but if the responder rejects, neither player receives anything. This game is often used to  
233 study fairness, social norms, and negotiation behavior (12).

234 In this game, ChatGPT is asked either to propose to divide \$100 between itself and a responder, or to respond to a proposal  
235 made by the other player. The difference between the Ultimatum and the Dictator game is that the responder no longer blindly  
236 accepts the proposal, and if they reject it, both players will get \$0 no matter what the original proposal was.

237 When playing the proposer alone, ChatGPT-3’s proposal of the amount given to the other player follows a distribution with  
238 mode \$40 ( $min = 30, max = 70, \mu = 45.2, \sigma = 12.1$ ). Compared with the Dictator game, ChatGPT-3 gives significantly more  
239 ( $p < 0.001$ ) to the responder when the other party’s decision jointly decides the outcome.

240 When specifically requested to provide an explanation of the decision, ChatGPT’s decision does not present a significant  
241 shift. The most common decision is still \$40 ( $min = 30, max = 70, \mu = 42.9, \sigma = 10.8, p = 0.54$ ). Like in the Dictator game,  
242 ChatGPT-4 makes a unanimous decision of an equal split (\$50,  $p = 0.001$ ), similar to that of ChatGPT Plus, and ChatGPT  
243 Free’s behavior is much more similar to that of ChatGPT-3 ( $mode = 40, min = 30, max = 70, \mu = 47.3, \sigma = 16.8, p = 0.905$ ).

244 Different context also has an effect on ChatGPT’s decision. When asked to explicitly explain the decision, ChatGPT-4  
245 presents a slight variance in its decision ( $\sigma = 2.6$ ). A similar variance is also observed when ChatGPT-4 is asked to play  
246 different roles: as a mathematician or an investment fund manager, it occasionally proposes a smaller split to the responder,  
247 and as a journalist, it occasionally proposes a greater split to the responder ( $max = 100, \sigma = 9.5$ ). None of these variations is  
248 statistically significant. The impact of roles to ChatGPT-3 is greater: the median of proposal decreases as a mathematician,  
249 centralizes as a public relation specialist, and increases as an investment fund manager.

250 When ChatGPT is asked to play the responder, it is asked the lowest amount that it is willing to accept. ChatGPT-3’s  
251 response still follows a bell-shaped distribution, with a mode of \$30 ( $\mu = 32.5, \sigma = 14.7$ ), a minimum of \$1, and a maximum of  
252 \$50. The mode, mean, min, and max are all lower than the statistics when it plays the proposer. A similar pattern is observed  
253 from ChatGPT Free ( $\mu = 30, \sigma = 11.9, p = 0.400$ ).

254 ChatGPT-4 presents a significantly different behavior, where the decisions follow a two-mode distribution that concentrates  
255 on two sides (\$1 and \$50). Requiring explanation does not change the two modes and only affects the distribution of the  
256 middle range. ChatGPT Plus presents an even more extreme trait, and in a dominating majority of sessions, it is willing to  
257 accept as low as \$1 as the responder ( $p \ll 0.001$ ). The choice of \$1 aligns well with the rational decision of the game. Indeed,  
258 some keywords that stand out from the responses include *nothing*, *better (than)*, and *something*.

259 The context of different occupations also affects ChatGPT’s decision in this game. As a mathematician, ChatGPT-4’s  
260 decision concentrates on the rational choice (\$1,  $p \ll 0.001$ ); as a public relation specialist, a journalist, an investment fund  
261 manager, or a legislator, its decision becomes less bipolarated, and the median shifts towards the right (favoring fairness more  
262 than rationality), although not statistically significant (Fig. S6c). Similar shifts are observed for ChatGPT-3 when inquired in  
263 the context of corresponding occupations.

264 The above experiment queries ChatGPT in independent sessions, so its response to one question is not interfered by its  
265 decision or memory about the other question. An interesting question is whether their exposure/response to one question  
266 (corresponding to one role as proposer or responder) affects their decision for a follow-up question that corresponds to the  
267 other role. We expose ChatGPT to both roles in the same session, asking it to respond to one question and then respond to  
268 the other question. We find that being exposed to one scenario does influence ChatGPT-3’s and ChatGPT-4’s responses in the  
269 other scenario, compared with the results obtained from independent sessions. Having been the proposer first (Fig. S7b), the  
270 distribution of ChatGPT-3’s response as the responder becomes narrower, with the \$50 group shifting left towards \$40. A  
271 similar shift presents in ChatGPT-4’s response, where the \$50 group moves towards \$20 and the median has reduced to \$15  
272 ( $p \ll 0.001$ ). Having been the responder first (Fig. S7a), ChatGPT-3’s response as the proposer also becomes more generous,  
273 while ChatGPT-4’s unanimous decision is not affected.

274 **D. The Trust Game.** The Trust Game is a two-player game that investigates trust and reciprocity. In this game, the first player  
275 (the investor– the trustor) is given a sum of money and can invest any amount in the second player (the banker– the trustee).  
276 The amount sent is multiplied by the game host, and the banker then decides how much, if any, to return to the investor. This  
277 game is used to study trust, reciprocity, and social norms (14).

278 In the game, ChatGPT is also asked to play one of two roles: as an investor or as a banker. As an investor, it is asked  
279 to decide how much (from \$0 to \$100) to invest in the banker (which is expected to generate a profit), who may return the  
280 entire revenue or nothing to the investor. ChatGPT-3’s decision follows a distribution that peaks at a moderate value of  
281 \$30 ( $min = 20, max = 60, \mu = 36.3, \sigma = 12.7$ ). When explicitly asked to explain the decision, its decision does not present a  
282 significant shift but shows a wider spread ( $min = 20, max = 100, \sigma = 16.0, p = 0.439$ ). Different occupation roles also show an  
283 effect, where ChatGPT-3 makes a larger investment as a mathematician, a public relation specialist, or a journalist.

284 ChatGPT-4 acts as if it has significantly more trust in the banker ( $p = 0.006$ ). Its decision follows a distribution peaked  
285 at investing half (\$50) of the endowment, which also appears to be the maximum investment it makes. It tends to invest  
286 slightly more when specifically required to provide an explanation. The distributions are also mildly affected by the assumed  
287 occupations, whereas for a mathematician or an investment fund manager, the distribution is more spread-out, and it even  
288 occasionally invests the entire endowment. The difference is not significant except for being a public relation specialist  
289 ( $p = 0.05$ ).

290 ChatGPT Plus invests less than ChatGPT-3, with a distribution peaked at \$10 ( $min = 0, max = 100, \mu = 19.6, \sigma =$   
291  $23.0, p \ll 0.001$ ). The keywords standing out from the responses include *risk*, *return*, *desires*, *losing*, *minimizes*, *guarantee*, and  
292 *control*. Once again, ChatGPT Free’s decision distribution is more similar to that of ChatGPT-3 ( $mode = 30, min = 10, max =$   
293  $60, \mu = 36.7, \sigma = 12.1, p = 0.792$ ).

294 In the second scenario, ChatGPT is asked to play the banker and decides what proportion of the total revenue is returned  
295 to the investor, which could range all the way from nothing (\$0) to the entire revenue (original investment + 2x profit). We  
296 find that the two common strategies that ChatGPT-4 uses are to return 1) the original investment plus half of the profit, or 2)  
297 half of the revenue. The former returns more to the investor than the latter. For ChatGPT-4, the former is more frequently  
298 used, and requiring the explanation makes this strategy even more dominating. Among the assumed occupations, ChatGPT-4  
299 returns more generously to the investor as a public relation specialist or an investment fund manager.

300 When the size of the investment increases, there is a shift in ChatGPT-4’s decision. When the investment was \$100 (out of  
301 \$100), the second and less generous strategy (returning half revenue) becomes the new dominant (Fig. S6e).

302 ChatGPT-3’s decision centers on the first strategy, with small probabilities of giving more generous (even the entire revenue)  
303 returns. This strategy is robust to the size of the investment. When the investment was \$10, \$50, and \$100, ChatGPT-3’s  
304 decision follows a distribution peaked at \$20 ( $\mu = 22.5, \sigma = 4.3$ ), \$100 ( $\mu = 112.5, \sigma = 27.5$ ), and \$200 ( $\mu = 239.2, \sigma = 40.7$ )  
305 (Fig. S6f). Requesting an explanation makes the distribution even more centralized, but the extremely generous behavior  
306 (returning the whole revenue to the investor) does not vanish. Being a public relation specialist, a journalist, or an investment  
307 fund manager, ChatGPT-3 tends to pay back more generously than in the default role.

308 Being exposed to a different role influences ChatGPT’s decision in the other role. When inquired about both scenarios  
309 (investor and banker) in the same session, ChatGPT’s responses for the second role deviate from the distribution observed in  
310 the independent sessions. In particular, after playing the role of the banker, both ChatGPT-3’s and ChatGPT-4’s investments  
311 increased. After playing the role of the investor, both ChatGPT-3’s and ChatGPT-4’s decisions as the banker become more  
312 generous.

313 **E. The Bomb Risk Game.** In the Bomb Risk Game, a player decides how many out of 100 boxes to collect, one of which contains  
 314 a ‘bomb.’ Earnings increase linearly with the number of boxes a player decides to open, but drop to zero if the ‘bomb’ is hit.  
 315 The game is designed to measure risk attitudes (15).

316 In this game, ChatGPT is asked to open a number of boxes out of 100, among which a bomb is randomly placed. If the box  
 317 that contains the bomb is not opened, then the player earns points that equal the number of boxes they open. If the box that  
 318 contains the bomb is opened, the bomb explodes and the player gets zero points. In every new round, a new set of 100 boxes is  
 319 provided and the bomb is placed at random. Opening 50 boxes is where the expectation of points gain is maximized, and the  
 320 decision at each round should be independent of the decisions/results in previous rounds.

321 When first exposed to the game, both ChatGPT-3’s and ChatGPT-4’s decisions peak at opening 50 boxes (ChatGPT-3:  
 322  $\mu = 39.8, \sigma = 16.3$ , ChatGPT-4:  $\mu = 57.6, \sigma = 17.4$ ), demonstrating a rational and risk neutral pattern of behavior. Despite  
 323 the mode, the distribution of ChatGPT-4 is significantly more risk-loving than that of ChatGPT-3 ( $p < 0.001$ ), with 13.8%  
 324 of the chance of opening a maximum number (99) of boxes. When playing the game for multiple rounds, its decision in the  
 325 following round is influenced by the outcome of the previous rounds. As long as the bomb does not explode in the previous  
 326 round, ChatGPT-3’s decision in the second round becomes even more concentrated on 50 boxes ( $\mu = 48.8, \sigma = 7.7$ ). When the  
 327 bomb does explode in the first round, ChatGPT-3 tends to be more conservative in the second round, with a new mode of  
 328 opening 25 boxes ( $\mu = 26.4, \sigma = 10.8$ ). If the bomb explodes again in the second round, it decides to open even fewer boxes  
 329 ( $mode = 10, \mu = 14.8, \sigma = 11.0$ ) in the third round. If the bomb exploded in the second round but not in the first round, the  
 330 average number of boxes opened in the third round drops from 50 to 43.1 even though the mode is still 50 ( $mode = 50, \sigma = 25.0$ ).

331 **F. Public Goods.** In the Public Goods game, each participant is given an initial endowment and can contribute any portion to  
 332 a public good project. The total amount raised for the project is then multiplied by a factor and distributed equally among all  
 333 participants, regardless of their individual contributions. This setup allows for the exploration of altruistic contributions, the  
 334 free-rider problem, and social dilemmas (16).

335 In the game, ChatGPT is asked to decide how much money from 0-\$20 to contribute to a public project as one of four  
 336 participants, the personal payoff of which is the sum of the amount not invested and 50% of the group contribution. It therefore  
 337 makes a profit when the contribution from the rest of the group is greater than its contribution. It is the most common for  
 338 both ChatGPT-3 and ChatGPT-4 to invest half of the endowment into the public goods project. When ChatGPT-3 receives a  
 339 greater payoff than the original endowment, it tends to increase its contribution in the next rounds despite whether the other  
 340 players made a larger or smaller contribution (or received a higher payoff). ChatGPT-4 on the other hand, makes a consistent  
 341 contribution in most sessions despite the contribution of other players. In sessions where its contributions do increase over  
 342 rounds, the increases tend to be smaller than those of ChatGPT-3.

343 **G. Prisoner’s Dilemma.** The Prisoner’s Dilemma is a fundamental game in game theory and behavioral economics, illustrating  
 344 the tension between individual and collective rationality. The game highlights that individuals, acting in their own self-interest,  
 345 can end up with worse outcomes than if they had cooperated (17–20).

346 In our version from MobLab, the framing is as follows. Two players are separately deciding whether to play the ‘push’  
 347 (cooperate) or ‘pull’ (defect) card. If both push, they collectively earn a higher payoff; if one pulls and the other pushes, the  
 348 defector gets all the payoff; if both pull, they both receive less payment. The payoff matrix is displayed in Table S3.

**Table S3. Payoff matrix of the Prisoner’s Dilemma game. The first numbers in cells are the payoffs of Player A, and the second numbers are the payoffs of Player B.**

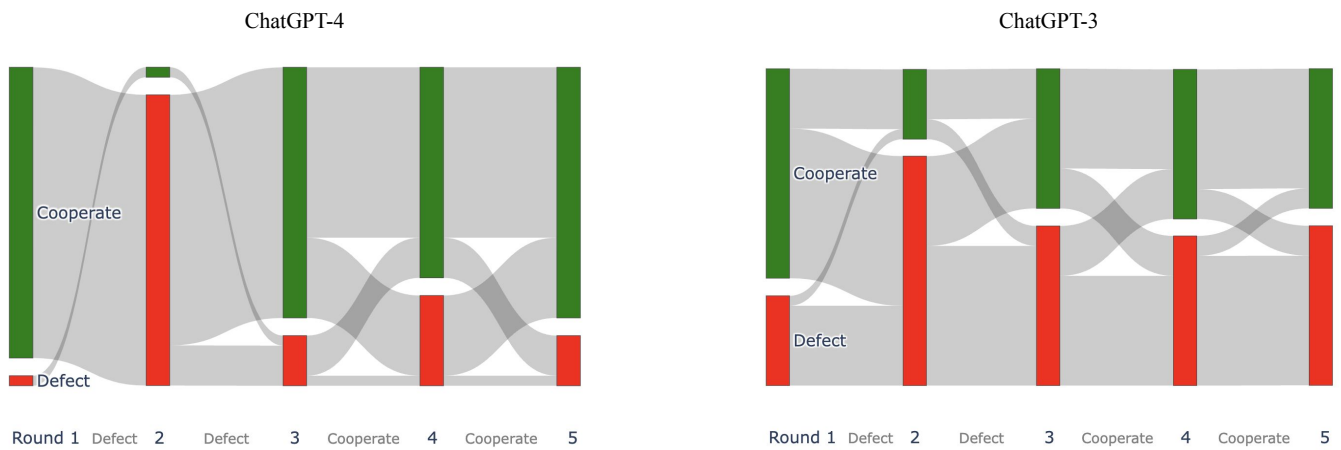
		Player B	
		Push	Pull
Player A	Push	\$400, \$400	\$0, \$700
	Pull	\$700, \$0	\$300, \$300

349 In the five-round game (Fig. S8, for the first round, ChatGPT attempts to cooperate with the other player in a majority  
 350 of the sessions (21/30 for ChatGPT-3, and 29/30 for ChatGPT-4). When the other player chooses to defect, however, the  
 351 majority of its decisions in the next round quickly turn into “Pull” (defect), consistent with the ‘tit for tat’ pattern. Such a  
 352 “punishing” strategy does not last if the other player continues to play “Pull.” Instead, the ratio of “Push” bounces back to  
 353 25/30 for ChatGPT-4 and 14/30 for ChatGPT-3, once again trying to incentivize the other player to coordinate.

354 When the other player plays two “Push” cards (chooses to cooperate) in a row in the third and fourth rounds, both ChatGPT-  
 355 3’s and ChatGPT-4’s decisions become relatively stable and end up cooperating in 14/30 of the instances (ChatGPT-3) and  
 356 25/30 of the instances (ChatGPT-4) after all five rounds.

## 357 References

358 1. Goldberg LR (1993) The structure of phenotypic personality traits. *American psychologist* 48(1):26.  
 359 2. McCrae RR, Costa Jr PT (1997) Personality trait structure as a human universal. *American psychologist* 52(5):509.  
 360 3. Roccas S, Sagiv L, Schwartz SH, Knafo A (2002) The big five personality factors and personal values. *Personality and*  
 361 *social psychology bulletin* 28(6):789–801.  
 362 4. McCrae RR, Costa Jr PT (2008) The five-factor theory of personality. *Theoretical Perspectives*.



**Fig. S8.** Five-round Prisoner's Dilemma Game. A large proportion of AI decisions switch from cooperation to defection if the other player defects in the first round. However, a significant portion reverts back to cooperation in the third round even if the other player continues to defect. The proportion of cooperation becomes relatively stable in the following rounds. We executed 5 rounds in total. Partner actions in each round (1-4) are observed after the player's action and recorded below each panel.



- 363 5. Goldberg LR (1992) The development of markers for the big-five factor structure. *Psychological assessment* 4(1):26.
- 364 6. Eloundou T, Manning S, Mishkin P, Rock D (2023) Gpts are gpts: An early look at the labor market impact potential of  
365 large language models. *arXiv preprint arXiv:2303.10130*.
- 366 7. Lin PH, et al. (2020) Evidence of general economic principles of bargaining and trade from 2,000 classroom experiments.  
367 *Nature Human Behaviour* 4(9):917–927.
- 368 8. OpenAI (2023) Gpt-4 technical report.
- 369 9. Ouyang L, et al. (2022) Training language models to follow instructions with human feedback. *Advances in Neural*  
370 *Information Processing Systems* 35:27730–27744.
- 371 10. McFadden D (1963) Constant elasticity of substitution production functions. *The Review of Economic Studies* 30(2):73–83.
- 372 11. McFadden D, , et al. (1973) Conditional logit analysis of qualitative choice behavior. *Frontiers in Econometrics*.
- 373 12. Güth W, Schmittberger R, Schwarze B (1982) An experimental analysis of ultimatum bargaining. *Journal of economic*  
374 *behavior & organization* 3(4):367–388.
- 375 13. Forsythe R, Horowitz JL, Savin NE, Sefton M (1994) Fairness in simple bargaining experiments. *Games and Economic*  
376 *behavior* 6(3):347–369.
- 377 14. Berg J, Dickhaut J, McCabe K (1995) Trust, reciprocity, and social history. *Games and economic behavior* 10(1):122–142.
- 378 15. Crosetto P, Filippin A (2013) The “bomb” risk elicitation task. *Journal of risk and uncertainty* 47:31–65.
- 379 16. Andreoni J (1995) Cooperation in public-goods experiments: kindness or confusion? *The American Economic Review* pp.  
380 891–904.
- 381 17. Von Neumann J, Morgenstern O (1947) *Theory of games and economic behavior*, 2nd rev. (Princeton university press).
- 382 18. Schelling TC (1958) The strategy of conflict. prospectus for a reorientation of game theory. *Journal of Conflict Resolution*  
383 2(3):203–264.
- 384 19. Rapoport A, Chammah AM (1965) *Prisoner’s dilemma: A study in conflict and cooperation*. (University of Michigan  
385 press) Vol. 165.
- 386 20. Andreoni J, Varian H (1999) Preplay contracting in the prisoners’ dilemma. *Proceedings of the National Academy of*  
387 *Sciences* 96(19):10933–10938.