

---

# Show, Don't Tell: Evaluating Large Language Models Beyond Textual Understanding with ChildPlay

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 The evaluation of Large Language Models (LLMs) often focuses on linguistic  
2 tasks, yet such assessments may not fully capture the models' general reasoning  
3 capabilities. We explore the hypothesis that LLMs, such as GPT-3.5 and GPT-4,  
4 possess broader cognitive functions, particularly in non-linguistic domains. Our  
5 approach extends beyond standard linguistic benchmarks by incorporating games  
6 like Tic-Tac-Toe, Connect Four, and Battleship, encoded via ASCII, to assess strategic  
7 thinking and decision-making. To evaluate the models' ability to generalize  
8 beyond their training data, we introduce two additional games. The first game,  
9 LEGO Connect Language (LCL), tests the models' capacity to understand spatial  
10 logic and follow assembly instructions. The second game, the game of shapes,  
11 challenges the models to identify shapes represented by 1s within a matrix of zeros,  
12 further testing their spatial reasoning skills. This "show, don't tell" strategy uses  
13 games to potentially reveal cognitive capabilities rather than simply querying the  
14 models. Our results indicate that despite their proficiency on standard benchmarks  
15 and temperature settings, GPT-3.5 and GPT-4's abilities to play and reason about  
16 fully observable games without pre-training is mediocre. Both models fail to  
17 anticipate losing moves in Tic-Tac-Toe and Connect Four, and they are unable to  
18 play Battleship correctly. While GPT-4 shows some success in the game of shapes,  
19 both models struggle with the assembly tasks presented in the LCL game. These  
20 results suggest that while LLMs like the GPT models can emulate conversational  
21 proficiency and basic rule comprehension, their performance in strategic gameplay  
22 and spatial reasoning tasks is limited in cognitive flexibility and generalization.  
23 Importantly, this reveals a blind spot in current LLM benchmarks that we highlight  
24 with our gameplay benchmark suite ChildPlay (GitHub Repository). Our findings  
25 provide a cautionary tale about claims of emergent intelligence and reasoning  
26 capabilities of LLMs that are roughly the size of GPT-3.5 and GPT-4.

## 27 1 Introduction

28 Typically, LLMs are transformer-based models that process input text and generate output text in a  
29 coherent and contextually appropriate manner. They utilize the self-attention mechanism to weigh  
30 the importance of different words in a sentence relative to each other [33, 6]. Input text is tokenized,  
31 converted into vectors using embeddings, and processed through transformer layers that calculate  
32 attention scores to dictate focus on relevant tokens [33, 6, 12]. The model then selects the next token  
33 based on learned distributions, iteratively generating an arbitrarily long sequence of text [33, 6, 12].  
34 With their enormous parameter counts, from Alpaca with 7 billion parameters [29], to LLaMA with  
35 65 billion [31] or even PaLM and its 540 billion parameters [11], these neural networks have learned  
36 to model complex linguistic abstractions, capturing patterns in syntax, semantics, pragmatics, and  
37 even elements of style and tone [6, 7, 21].

38 Benchmarks for evaluating Large Language Models (LLMs) have been designed to assess compre-  
39 hension, generation, and adaptability across a wide range of language tasks. Datasets like SQuAD,  
40 GLUE, BIG-bench, and the lm-evaluation-harness offer various test types, including multiple-choice  
41 questions, reading comprehension exercises, and dialogue completion tasks. These benchmarks  
42 deploy metrics such as response correctness, language generation fluency, and the ability to maintain  
43 contextually relevant dialogue [22, 34, 2, 26]. Other benchmarks like SuperGLUE, ANLI, Truth-  
44 fulQA, and HellaSwag have been developed to evaluate different aspects of LLM performance, such  
45 as natural language understanding, commonsense reasoning, and factual knowledge about diverse  
46 topics [35, 20, 18, 37].

47 Recent studies have explored alternative approaches to evaluate LLMs' reasoning abilities in non-  
48 linguistic modalities. Liga and Pasetto modeled the game Tic-Tac-Toe using ASCII characters, pitting  
49 LLMs against the minimax algorithm to observe emergent features, which, according to the authors,  
50 might be akin to consciousness. The minimax algorithm is widely considered the optimal algorithm  
51 for playing tic-tac-toe, as it guarantees a win or draw against a perfect opponent [27, 1]. While LLMs  
52 performed well in some instances, they generally failed to win against the minimax algorithm, often  
53 resulting in a draw [17]. Topsakal and Harper [30] used Tic-Tac-Toe encoded with list and illustration  
54 prompts in their study. They found that while GPT-4 secured the most wins, it did not always win,  
55 indicating that GPT models cannot play Tic-Tac-Toe optimally. This contradiction raises the question:  
56 can we truly say the model knows how to play Tic-Tac-Toe if it can explain optimal strategies (see  
57 Appendix A.3) but does not consistently win? Or is its performance merely the result of probabilistic  
58 outcomes?

59 Some critical studies have highlighted the need for caution in interpreting LLMs' capabilities through  
60 benchmarking. Lappin et al. assessed their strengths and weaknesses, finding that they excel at  
61 many language tasks but struggle with deeper reasoning, world knowledge integration, and context  
62 understanding beyond local co-occurrences [16]. And Zečević et al. argued that LLMs may discuss  
63 causality but lack true causal reasoning based on interventions and counterfactuals [38].

64 Bender et al. argue that the lack of transparency and potential risks associated with these large,  
65 opaque models raise concerns about their trustworthiness and accountability [3]. While the criticism  
66 of Bender et al. focuses on the social dimension of the problem of interpretability and trustworthiness,  
67 recent work by Schaeffer et al. criticizes emergent capabilities and the perceived intelligence of LLMs.  
68 They suggest that some claimed "emergent abilities" of LLMs may be an artifact of the choice  
69 of evaluation metric, rather than fundamental changes in model behavior [23]. Their analyses  
70 demonstrate how the use of nonlinear or discontinuous evaluation metrics can create the illusion of  
71 emergent abilities, even when the underlying model performance changes smoothly and predictably  
72 with scale.

73 This critique of the evaluation metrics used in assessing LLMs invites a deeper exploration of general  
74 intelligence - specifically how it can be reliably measured and observed in AI through rigorous  
75 and realistic tests that extend beyond linguistic prowess to include broader cognitive functions. If  
76 we must define general intelligence (GI), one is to use the "g factor," which refers to the ability to  
77 reason, plan, solve problems, think abstractly, and learn quickly across a wide range of domains  
78 [24, 4, 36, 9, 8]. GI then involves higher-order cognitive processes that go beyond specific skills or  
79 knowledge domains [14, 15].

80 A critical issue that arises in analysing the reasoning capabilities of large and opaque models like the  
81 GPT series, is training-test set cross-contamination, which becomes increasingly problematic for the  
82 most advanced models [6]. The massive training datasets used, comprising extensive portions of the  
83 internet, are often untraceable and completely anonymous to researchers outside the initial developer  
84 groups, to some extent even to the developers themselves, making replication studies impossible  
85 [6, 13]. The exact amount and identity of data used to train models like GPT-3.5 or GPT-4 has not  
86 been publicly disclosed, posing a risk of rendering current benchmarking efforts meaningless due to  
87 cross-contamination.

88 Researchers have attempted to counter the contamination problem using N-Gram Overlap as a metric  
89 for detection, by eliminating or withholding results for tests where answers were present in the  
90 training data [6]. However, this method has been criticized. Blodgett et al. point out, for example,  
91 that such heuristic approaches to mitigating biases in NLP systems can be problematic and may not  
92 fully address the underlying challenges [5]. The method is also limited in that it fails to consider  
93 the context in which N-Grams appear and may discount synonymous or analogous text worded

94 differently. Additionally, the decision to use a 200-character window around detected N-Grams is  
 95 arbitrary and may not accurately reflect the influence of surrounding text on model learning.

96 In this work we introduce ChildPlay, a suite of non-language-based games like Tic-Tac-Toe, Connect-  
 97 Four, Battleship, LEGO Connect Language, and the game of Shapes, to assess reasoning, strategic  
 98 capabilities, symbolic reasoning, and pattern recognition abilities of large language models (LLMs)  
 99 beyond traditional linguistic modalities. Games provide structured environments with clear success  
 100 criteria, making them suitable for evaluating strategic thinking, planning, and long-term decision-  
 101 making of LLMs [25, 17, 30]. Their dynamic and adversarial nature resembles real-world scenarios,  
 102 assessing generalized intelligence and reasoning beyond the training domain [25, 17, 30]. We encode  
 103 these games using ASCII representations to minimize dataset contamination issues prevalent in  
 104 contemporary LLM benchmarks [6, 17].

## 105 2 Experiments

106 Specific tasks in the BIG-bench benchmark [2], among others, are categorized as either zero-shot,  
 107 one-shot, or multi-shot [6]. Our tasks fit the zero-shot category, as models are given only a brief  
 108 explanation at inference time with no examples for playing beyond the explained formalism. To  
 109 demonstrate the reasoning capabilities of LLMs beyond their training data, we focus on a modality not  
 110 explicitly trained for: spatial reasoning about ASCII sequences. An agent capable of true abstraction  
 111 should be able to encode and interpret these sequences if the rules are explained or known.

112 For this purpose, we developed several tasks, including LEGO assembly, ASCII games of Tic-Tac-  
 113 Toe, Connect-Four, and Battleship, as well as identifying simple geometrical shapes represented as 1s  
 114 in 15-sided grids of 0s. The same models were deployed over all experiments, namely *gpt-3.5-turbo-*  
 115 *1106*, and *gpt-4-1106-preview*, which in this paper are referred to as GPT-3.5 and GPT-4, respectively.  
 116 Every experiment was tested across different temperature settings ( $t$ ) per model, namely  $t=0$ ,  $t=0.5$ ,  
 117  $t=1$ , and  $t=1.5$ . When asked about their understanding of the tasks, GPT-3.5 and GPT-4 were able to  
 118 generate board states and explain the queried games, including their rules and optimal play. Thus, we  
 119 consider the tests valid: if the models are truly capable of reasoning, they should be able to play these  
 120 games optimally given that they "know" and are capable of explaining what playing optimally means  
 121 (see Appendix A.3). Experiments ran over night, at minimum taking a couple of minutes and at most  
 122 taking a few hours.

123 **LEGO Connect Language (LCL)** We invented a formal language we call LEGO Connect Language  
 124 (LCL). More specifically, we propose  $LCL_2$  as a language to instruct assembly in 2D on the  $x$  and  $y$   
 125 axis (this can easily be generalised to  $LCL_3$  - instructions along the  $x$ ,  $y$ , and  $z$  axis). The models  
 126 were given instructions and their output was fed through a visualizer script to generate the images  
 127 contained in this work. Only  $2 \times 4$  pieces were allowed. A piece  $P$  (see Fig 1) is then defined as a  
 128 tuple  $P = (l, w, (x, y), c, h)$ . A construction,  $M$ , is then a valid construction in  $LCL_2$  if no pieces  
 129 are overlapping and all pieces are connected to other pieces. Namely, a Lego piece is connected  
 130 through interlocking pegs, not by merely touching sides. And secondly, two Lego pieces overlap  
 131 when they share the same  $y$ -coordinate and any part of their length has the same  $x$ -coordinate.

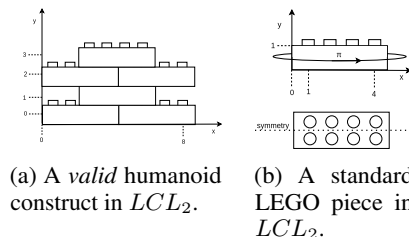


Figure 1: Introducing  $LCL_2$ .

132 **Game 1: Validity Testing** In this experiment, we evaluate the ability of different models to validate  
 133 the correctness of a given Lego construct. The constructs are generated to be either valid or invalid.  
 134 A construct is considered valid if there is no horizontal overlap between pieces, and pieces must  
 135 connect via overlapping pegs such that the whole assembly is connected (no floating pieces). The

136 models, namely GPT-4 and GPT-3.5, are then tasked with predicting the validity of each construct.  
 137 The evaluation metric for this experiment was the proportion of correct validations, measured across  
 138 different temperature settings.

139 **Game 2: Construct Generation** In this experiment, the models attempt to generate valid LCL  
 140 constructs. Each construct description consists of a list of tuples, where each tuple specifies the  
 141 coordinates and color of a Lego piece. The models generated these constructs based on prompts and  
 142 the validity of the constructs was automatically evaluated. The metric for this experiment was the  
 143 proportion of valid constructs generated, measured across different temperature settings.

144 We automatically produced 800 images for the validity test, half valid and half invalid ones. Then  
 145 each model was queried to produce 100 images at each temperature setting, which we then checked  
 146 for validity. We believe our use of LCL is related to the tests found in Bubeck et al. [7], where  
 147 JavaScript or LaTeX was used to prompt GPT-4 to produce images. However, while the images in  
 148 Bubeck et al. [7] included common examples such as letters, a car, a truck, a cat, a dog, a person,  
 149 a pig, a house, and a unicorn, all of which are likely represented in the training data in JavaScript  
 150 or LaTeX, LCL challenges the model to step outside of its learned data distributions by remaining  
 151 abstract.

152 **Three Board Games: Tic-tac-toe, Connect-four, and Battleship** In the case of the three board  
 153 games, each new board state was accompanied by the introductory game explanation sent through the  
 154 OpenAI API in a zero-context testing environment. The models were provided with the current board  
 155 state and an opponent making moves at random, with the LLM always playing as the first player,  
 156 which is advantageous in all three games. Context beyond the initial instruction and the current  
 157 board state was deemed irrelevant since these games are fully observable, meaning every board state  
 158 contains all the necessary information to play optimally. The input to the game was simply two  
 159 scalars for the row-column pair or just a scalar for the column number in the case of connect-four.

160 For the battleship game, ships ('S') were randomly initialized, always horizontally, with varying sizes  
 161 spanning between 2 and 5 cells. When there is a hit by either player, the position is marked with an  
 162 'X' on both players' boards. If the guess was a miss, an 'O' is placed on the position instead.

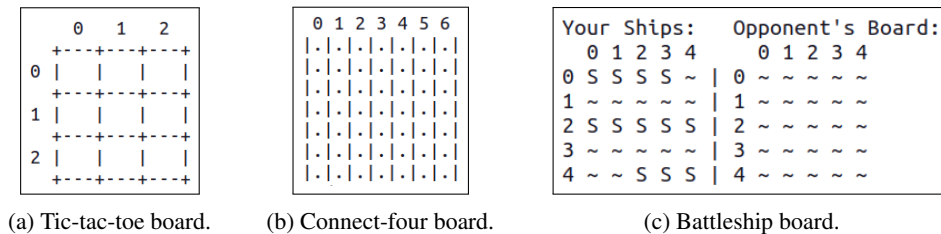


Figure 2: Initial board states as presented to the LLM (the ship positions in the Battleship board are randomised with every initialisation, including ship length).

163 **The Game of Shapes** In the case of the game of shapes, preliminary work involved probing the  
 164 models to determine what geometric shapes they consider basic by prompting them multiple times.  
 165 The first three shapes consistently mentioned were square, circle, and triangle (not necessarily in that  
 166 order). The game then consists of finding a basic geometric shape "hidden" behind 1s within a matrix  
 167 of 0s in a multiple-choice fashion. Four shapes were used as options: the circle, the rectangle, the  
 168 triangle, and the cross, but only the latter three were ever shown to the model (cf. Fig. 3).

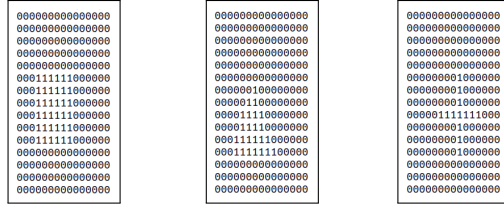


Figure 3: Matrices containing shapes used in the game of Shapes.

169 **3 Results**

170 As previously stated, Tic-Tac-Toe as a benchmark has been tackled before [17, 30]. Since it is quite  
 171 popular, we decided to replicate it before creating new games. But this time using an ASCII  
 172 encoding instead of a list of moves such that we can gauge spatial reasoning through symbolic  
 173 reasoning. For comparison with the model’s performance, Fig. 4 presents the Tic-Tac-Toe match  
 174 results of the *minimax* algorithm against the same random player the models played against. This  
 175 outcome creates a baseline for optimal play against a random player.

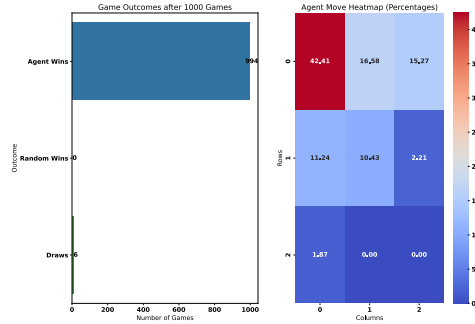


Figure 4: Minimax vs random player.

177 **Tic-tac-toe, Connect-four, and Battleship** To check for a win, we determine if the player  
 178 has successfully connected the winning number of pieces in a row on the board, which could be  
 179 horizontally, vertically, or diagonally. To detect missed and blocking moves, we simulate all potential  
 180 moves for the player by checking if placing a piece in any column leads to a win. If such a move  
 181 is found, and the player does not execute it on their turn, it is recorded as a missed win, if such a  
 182 move is found for the opponent and the player does not block it, we register it as missed blocking  
 183 move. We define *incorrect moves* to mean a move that was illegal, such as playing a position that has  
 184 already been played. This results in an immediate loss.

185 Fig. 5 encompasses comparative results from playing Connect-Four, Tic-Tac-Toe, and Battleship.  
 186 Each subfigure, 5a, 5b, and 14, respectively, outlines the number of games won by the models.

187 Unfortunately, the models were incapable of following the rules for the Battleship game, that is,  
 188 regardless of temperature, the models lose the large majority of games, with GPT-4 not winning a  
 189 single game due to incorrect moves (cf. Fig. 16). GPT-3.5 wins around 10% of the matches at low  
 190 temperatures, but none at higher temperatures, we refer to Fig. 14 in the Appendix A.1.3 instead.

191 It is notable that both GPT-3.5 and GPT-4 exhibit their poorest performance in both Connect-Four  
 192 and Tic-Tac-Toe at a temperature setting of 0, indicative of deterministic play that reflects the models’  
 193 learned strategies (Appendix A.1). The Random Player’s normal distribution across columns (Fig.  
 194 12) suggests a lower likelihood of countering GPT’s central strategies, in both games, but particularly  
 195 at Tic-Tac-Toe where GPT-3.5 commits more errors than GPT-4, significantly impacting outcomes  
 196 due to incorrect moves (Fig. 5b). These errors generally increase with temperature, probably due  
 197 to enhanced choice randomness (Fig. 10). This explains the lack of direct model losses from final  
 198 defeating moves since losses often result from illegal moves.

199 Average game moves, missed wins, and blocks in both Tic-Tac-Toe and Connect-Four are further  
 200 illustrated in Figs. 6a and 6b, highlighting a decrease in these metrics as temperature rises, suggesting  
 201 that higher settings potentially broaden the explored moves within the models' strategies. Conclu-  
 202 sively, neither model plays the games optimally, as evidenced by the considerable number of missed  
 203 wins and blocks. Both subfigures demonstrate that, as temperature increases, the number of missed  
 204 wins and blocks decreases. This might suggest that higher temperature settings potentially increase  
 205 the explored moves in the models' learned strategy, in case there is any. We can conclude the same as  
 206 before, namely that neither model can play Tic-Tac-Toe optimally given the number of missed wins  
 207 and missed blocks.

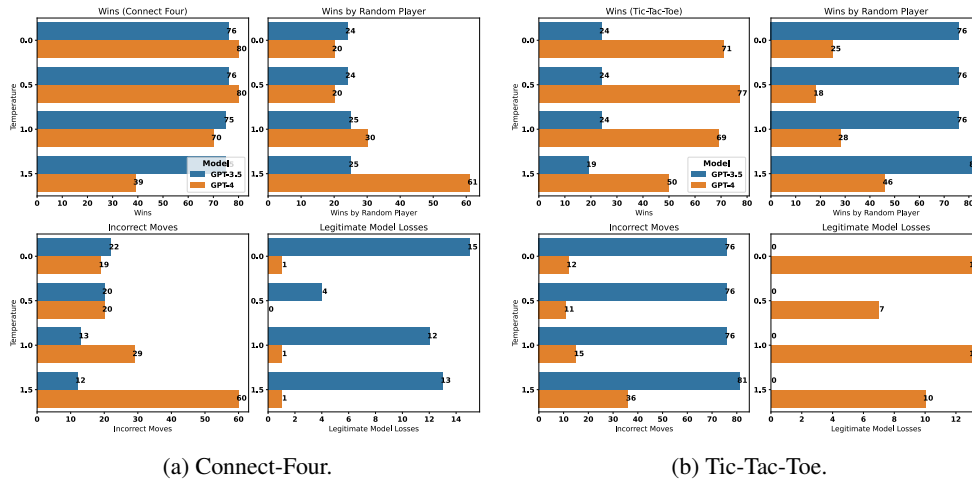


Figure 5: Incorrect Moves, Wins, and Losses Per Player in the Three Board Games.

208 The number of moves of GPT-3.5 and GPT-4 per game can be thought of as a measurement of stability  
 209 in gameplay, not just against the random player, but in general, given that a longer game entails that  
 210 the model is not losing to illegal moves or to its opponent. It increases linearly with temperature,  
 211 inversely correlated with performance measured by the decrease in missed wins and blocks. Tic-  
 212 Tac-Toe shows a linear improvement, whereas Connect-Four experiences an exponential boost in  
 213 performance from temperature 0 to 0.5, followed by a linear decline. The random player consistently  
 214 performs better against GPT-3.5 in Tic-Tac-Toe but loses more frequently in Connect-Four. Both  
 215 models struggle with blocking or seizing winning moves from the random player. An analysis of the  
 216 move heatmaps (cf. Appendix A.1) explains why winning Connect-Four against a random player  
 217 appears straightforward. As the model consistently places pieces in the same column, the probability  
 218 of the random player losing increases with the board size. However, even under these challenging  
 219 conditions, the random player still secures wins in at least 20% of the games played against GPT-4.

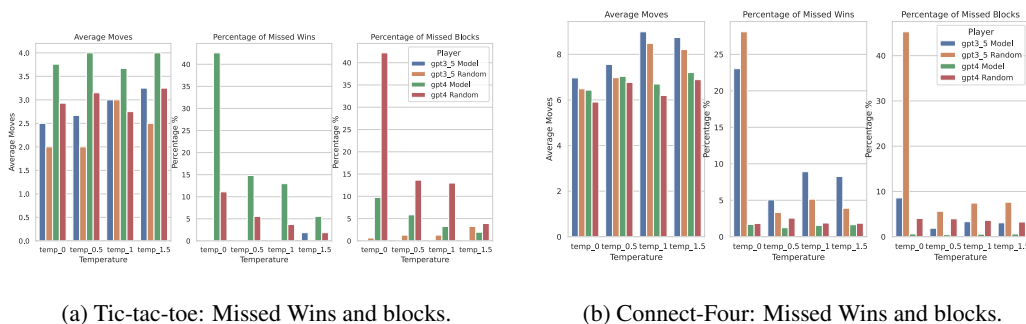
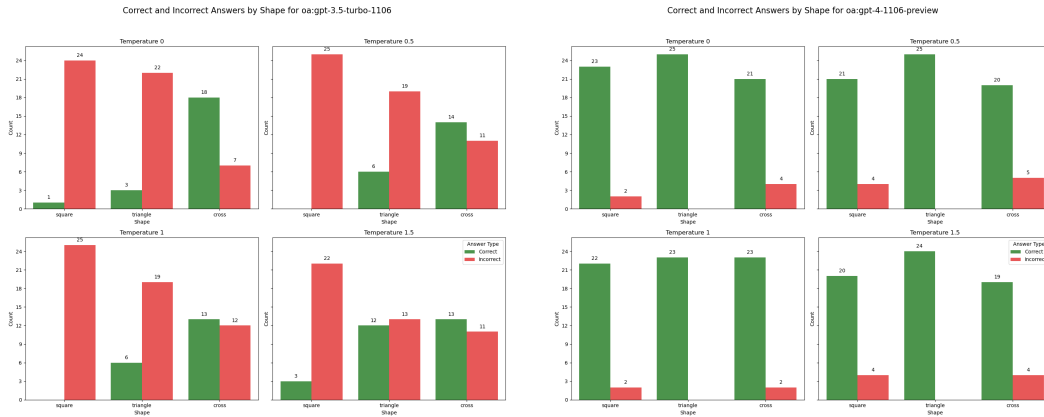


Figure 6: Average Moves, missed wins, and missed blocks for Tic-tac-toe and Connect-Four.

220 **Shapes** In the game of Shapes, a correct detection happens when the player’s selected shape  
 221 corresponds with the shape shown on the board. Players have four choices: "circle," "triangle,"  
 222 "square," and "cross". Notably, a circle is never actually displayed to the model, and the positions of  
 223 these choices are not randomized to test if the model displays any inherent bias for the question order.  
 224 This does not affect the outcome, since the game does not change across different sessions as it is  
 225 designed to operate within a single question-response framework.

226 In the shape detection tests, GPT-3.5’s performance was approximately equivalent to random chance  
 227 when identifying triangles and crosses, yet it completely failed to recognize squares. In contrast,  
 228 GPT-4 performed remarkably well, successfully identifying shapes with an accuracy of 80% or higher,  
 229 particularly proficient at recognizing triangles<sup>1</sup>.



(a) Results for the Shapes game, as played by GPT-3.5. (b) Results for the Shapes game, as played by GPT-4.

Figure 7: Experiment results for the Shapes game, comparing GPT-3.5 and GPT-4.

230 **LCL** In the game of LCL, both models systematically failed to respect the two rules, namely  
 231 that Lego pieces must be connected through interlocking pegs, not by merely touching sides, and  
 232 secondly, that no Lego pieces may overlap, which occurs when they share the same y-coordinate and  
 233 any part of their length has the same x-coordinate. For example, Figs. 8, 8a, and 8b show valid LCL  
 234 assemblies, while Figs. 8c and 8d show invalid LCL structures. Figs. 8a and 8b show valid LCL  
 235 assemblies, while subfigs. 8e and 8g show invalid output from GPT-3.5 generated at temperature 0.  
 236 While Fig.8f shows a valid output from GPT-4 at temperature 1.5. Other images (Figs. 8i, 8j, 8k, and  
 237 8l) are of invalid output<sup>2</sup>.

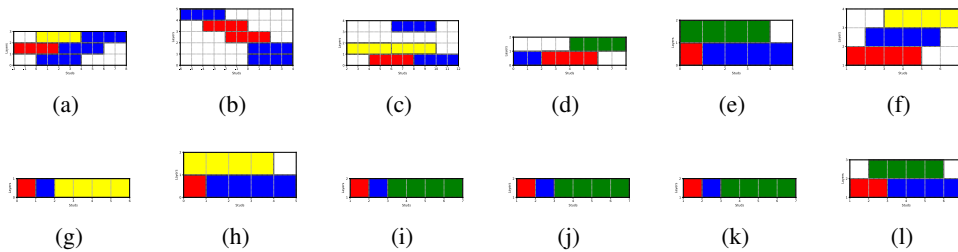


Figure 8: Structures automatically generated for the LCL validity test and structures generated by GPT-3.5 and GPT-4 for the construction generation test.<sup>3</sup>

<sup>1</sup>At higher temperatures, some of GPT-4’s responses were discarded by our parser when the model generated invalid Unicode output, and thus were not included in the final evaluation. This discrepancy is evident in Fig. 7b, for instance, where the sum of correct and incorrect choices does not total 25 at temperatures 1 and 1.5.

<sup>2</sup>Fig. 8i = GPT-4 at temperature 0, Fig. 8j = GPT-4 at temperature 0.5, Fig. 8k = GPT-4 at temperature 1, and Fig. 8l = GPT-4 at temperature 1.5.

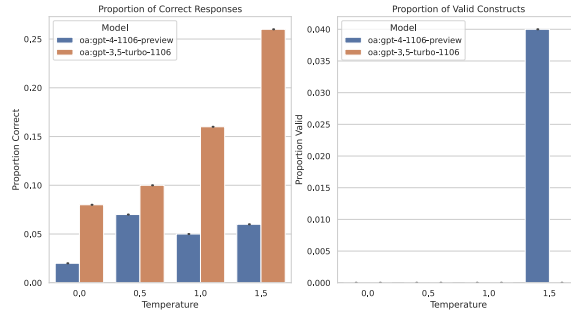


Figure 9: LCL results after 100 runs with 50/50 valid/invalid examples for the validity test and 100 experiments per temperature per model for the construction modality using 3 pieces.

238 Fig. 9 shows a roughly linear increase in the proportion of correct answers during the validity test as  
 239 a function of temperature. However, only GPT-4 produced a small minority of valid LCL constructs  
 240 (namely 0.04 of a total of 400 = 16), while GPT-3.5 did not manage to produce a single valid LCL  
 241 construct.

## 242 4 Discussion

243 In Tic-Tac-Toe, both models underperform relative to the minimax algorithm baseline, while showing  
 244 mixed performance at Connect-Four. GPT-4 performs unexpectedly well at the Shapes game, but  
 245 GPT-3.5 does very poorly. Also unexpectedly, both models fail to assemble or detect valid Lego  
 246 structures in the LCL game. In Battleship, the models' failure to follow game rules, especially at  
 247 higher temperature settings, indicates a significant limitation in their ability to understand and apply  
 248 structured game rules. The linear increase in the number of moves with temperature suggests that  
 249 higher temperatures lead to greater exploration of possible moves, but do not improve strategic  
 250 performance. The increase in missed wins and blocks with temperature further supports this, as  
 251 greater randomness in decision-making does not enhance the models' strategic play.

252 Overall, these results show that while GPT-3.5 and GPT-4 can play simple games to some extent, they  
 253 struggle with more complex tasks and do not consistently apply optimal strategies. The performance  
 254 gap between the models and the minimax algorithm highlights the limitations of current language  
 255 models in tasks requiring precise strategic reasoning and the failure to play Battleship and LCL  
 256 demonstrates a failure in rule adherence.

257 The primary aim of contemporary benchmarks for LLMs has been to assess these models through  
 258 adaptations of Turing's test [32], evaluating their capability to process and respond to language inputs  
 259 comparably to humans. However, defining the language problem solely in these terms may overlook  
 260 deeper complexities. While the transformer architecture in deep neural networks has enabled models  
 261 smaller than GPT-4 to exhibit what Wilhelm von Humboldt described as the "infinite use of finite  
 262 means" [19] or their ability to generate a potentially unlimited number of contextually relevant  
 263 sentences [28] (an idea popularised by Chomsky [10]), this does not necessarily imply that these  
 264 models have mastered a form of reasoning. Rather, they may simply be engaging in an advanced  
 265 form of pattern imitation.

### 266 4.1 Limitations and Future Work

267 Our proposed benchmark, ChildPlay, primarily uses binary (win/loss) outcomes for games, which  
 268 can be considered discontinuous metrics. Mathematically, these are expressed as:

$$\text{Metric}(x) = \begin{cases} 1 & \text{if win} \\ 0 & \text{if loss} \end{cases}$$

<sup>3</sup>Images in Fig. 8 were not directly produced by the GPT models. Instead, the formal descriptions of these images were generated by the models and subsequently passed to a script for rendering available in the GitHub Repository.



269 This formulation may exaggerate perceived capabilities by registering a full loss even if the model’s  
270 failure was marginal. We try to avoid this simplistic classification by registering, for example, the  
271 choice of moves on the board games (see Appendix A.1) as well as the count of missed blocks  
272 and missed wins (cf. Fig. 6). In contrast, tasks involving shape recognition or LCL could utilize  
273 more continuous metrics, providing a smoother performance gradient and potentially more accurate  
274 reflections of a model’s reasoning abilities.

275 Using discontinuous metrics in strategic games could manifest as sharp transitions in model evalua-  
276 tion:

$$\text{Performance}(N) = \delta(\text{outcome}_N - \text{threshold})$$

277 where  $\delta$  is the Dirac delta function, accentuating a sudden jump in perceived ability when the model  
278 first succeeds. Nonlinear metrics in the shape game or LCL tasks may not exhibit such abrupt  
279 transitions but could still misrepresent gradual improvements:

$$\text{Performance}(N) \approx \exp(-\alpha N^\beta)$$

280 where  $\alpha > 0$  and  $\beta < 0$  dictate the rate of improvement. This expression reflects smoother but  
281 potentially misleadingly slow progress.

282 Based on the perspective from Schaeffer et al. [23], one could argue that the games proposed in  
283 ChildPlay may not entirely reflect true generalization or emergent abilities. If these benchmarks are  
284 akin to nonlinear or discontinuous metrics, they might exaggerate the weaknesses or strengths of  
285 LLMs in strategic games. For instance, a sharp failure in a game like Tic-Tac-Toe might not mean the  
286 model lacks strategic reasoning universally but that it fails under the specific discontinuous conditions  
287 of the game setup, or of temperature. Such an assessment could lead to the erroneous conclusion that  
288 LLMs are generally poor at strategic decision-making when, in fact, they might only be unsuited to  
289 the specific scenarios or metrics used in ChildPlay.

290 Conversely, unlike continuous metrics that might smooth over deficiencies and give a misleading  
291 picture of gradual improvement, the use of clear, structured games as benchmarks could provide a  
292 direct assessment of an LLM’s cognitive and strategic abilities regardless of metric continuity. That  
293 is, given that the model has not been overfitted on the game.

## 294 5 Conclusions

295 Non-language-based tasks are important as they challenge models to demonstrate generalization  
296 across different information encodings or forms of input, and, most importantly, to actually delve  
297 into out-of-training-distribution topologies. Testing LLMs like GPT-4 (according to OpenAI, the  
298 current contender to AGI [7]) beyond the text they were primarily trained on via our "show, don’t  
299 tell" strategy, we demonstrate that it is still mediocre at best at very simple reasoning tasks that are  
300 outside of its training data. The models fail to play optimally at very simple games, such as tic-tac-toe,  
301 battleship, and connect-four. We also experimented with LEGO assembly, finding the LLMs still  
302 performing poorly. Mixed results were found at the task of interpreting geometric shapes from binary  
303 grids. These tasks are then designed to test reasoning without relying on language skills, such that  
304 the model cannot get by through parroting - it must be capable of playing the game. In the context of  
305 BIG-bench, our tasks would fit in the "non-language" category. Currently, this category shows 16  
306 active tasks, including some explicit ASCII recognition tasks, chess, and Sudoku, however, to the  
307 best of our knowledge, no task like ours [2]. Hence, we believe that ChildPlay is a useful addition to  
308 the suite of current established LLM benchmarks.

309 In general, this work is relevant in that developing games allows us to critically examine claims  
310 regarding a models’ ability to perform reasoning and problem solving regardless of the persistent  
311 problem of data contamination. In other words, we explore what the model knows by making it  
312 play games instead of asking it how to play them. Our results suggest that current LLMs show  
313 disappointing performance in terms of problem solving capabilities and reveal important aspects to  
314 be considered for future improvements.

## 315 References

- 316 [1] Shahd H. Alkaraz, Essam El-Seidy, and Neveen S. Morcos. Tic-tac-toe: Understanding the mini-  
317 max algorithm. 2020. URL <https://api.semanticscholar.org/CorpusID:218798654>.
- 318 [2] BIG bench authors. Beyond the imitation game: Quantifying and extrapolating the capabilities  
319 of language models. *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856.  
320 URL <https://openreview.net/forum?id=uyTL5Bvosj>.
- 321 [3] Emily M Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. On  
322 the dangers of stochastic parrots: Can language models be too big? *Proceedings of the 2021*  
323 *ACM Conference on Fairness, Accountability, and Transparency*, 2021.
- 324 [4] A. Binet and T. Simon. *The development of intelligence in children*. Baltimore: Williams &  
325 Wilkins, 1916.
- 326 [5] Su Lin Blodgett, Solon Barocas, Hal Daumé III, and Hanna Wallach. Language (technology)  
327 is power: A critical survey of "bias" in nlp. In *Proceedings of the 58th Annual Meeting of the*  
328 *Association for Computational Linguistics*, 2020.
- 329 [6] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal,  
330 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel  
331 Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M.  
332 Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz  
333 Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec  
334 Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. *CoRR*,  
335 abs/2005.14165, 2020. URL <https://arxiv.org/abs/2005.14165>.
- 336 [7] Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, John A. Gehrke, Eric Horvitz, Ece  
337 Kamar, Peter Lee, Yin Tat Lee, Yuan-Fang Li, Scott M. Lundberg, Harsha Nori, Hamid Palangi,  
338 Marco Tulio Ribeiro, and Yi Zhang. Sparks of artificial general intelligence: Early experiments  
339 with gpt-4. *ArXiv*, abs/2303.12712, 2023. URL [https://api.semanticscholar.org/](https://api.semanticscholar.org/CorpusID:257663729)  
340 [CorpusID:257663729](https://api.semanticscholar.org/CorpusID:257663729).
- 341 [8] J.B. Carroll. *Human cognitive abilities: A survey of factor-analytic studies*. New York:  
342 Cambridge University Press, 1993.
- 343 [9] R.B. Cattell. Theory of fluid and crystallized intelligence: A critical experiment. *Journal of*  
344 *Educational Psychology*, 54(1):1–22, 1963.
- 345 [10] Noam Chomsky. *Syntactic Structures*. Mouton and Co., The Hague, 1957.
- 346 [11] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam  
347 Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker  
348 Schuh, Kensen Shi, Sashank Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes,  
349 Yi Tay, Noam Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson,  
350 Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin,  
351 Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier  
352 Garcia, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David  
353 Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani  
354 Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayanan Pillai, Marie Pellat,  
355 Aitor Lewkowycz, Erica Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei  
356 Zhou, Xuezhi Wang, Brennan Saeta, Mark Diaz, Orhan Firat, Michele Catasta, Jason Wei,  
357 Kathy Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. Palm: scaling  
358 language modeling with pathways. *J. Mach. Learn. Res.*, 24(1), mar 2024. ISSN 1532-4435.
- 359 [12] John Fields, Kevin Chovanec, and Praveen Madiraju. A survey of text classification with  
360 transformers: How wide? how large? how long? how accurate? how expensive? how safe?  
361 *IEEE Access*, 12:6518–6531, 2024. URL [https://api.semanticscholar.org/CorpusID:](https://api.semanticscholar.org/CorpusID:266824505)  
362 [266824505](https://api.semanticscholar.org/CorpusID:266824505).
- 363 [13] L. Floridi and Massimo Chiriatti. Gpt-3: Its nature, scope, limits, and consequences. *Minds*  
364 *and Machines*, 30:681 – 694, 2020. URL [https://api.semanticscholar.org/CorpusID:](https://api.semanticscholar.org/CorpusID:228954221)  
365 [228954221](https://api.semanticscholar.org/CorpusID:228954221).

- 366 [14] Linda S Gottfredson. Why g matters: The complexity of everyday life. *Intelligence*, 24(1):  
367 79–132, 1997.
- 368 [15] A.R. Jensen. *The g factor: The science of mental ability*. Westport, CT: Praeger, 1998.
- 369 [16] Shalom Lappin. Assessing the strengths and weaknesses of large language models. *Unpublished*  
370 *Manuscript*, 2023.
- 371 [17] Davide Liga and Luca Pasetto. Testing spatial reasoning of large language models: the case of  
372 tic-tac-toe. *Unpublished Manuscript*, 2023.
- 373 [18] Bill Yuchen Lin, Maarten Sap, Ari Holtzman, Antoine Bosselut, Hannah Rashkin, and Yejin  
374 Choi. Truthfulqa: Measuring how models mimic human falsehoods. *Proceedings of the 60th*  
375 *Annual Meeting of the Association for Computational Linguistics*, 2022.
- 376 [19] William Merrill. Formal languages and neural models for learning on sequences. In *International*  
377 *Conference on Graphics and Interaction*, 2023. URL [https://api.semanticscholar.org/](https://api.semanticscholar.org/CorpusID:261101973)  
378 [CorpusID:261101973](https://api.semanticscholar.org/CorpusID:261101973).
- 379 [20] Yixin Nie, Adina Williams, Emily Dinan, Mohit Bansal, Jason Weston, and Douwe Kiela.  
380 Adversarial nli: A new benchmark for natural language understanding. *Proceedings of the 58th*  
381 *Annual Meeting of the Association for Computational Linguistics*, 2020.
- 382 [21] Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin,  
383 Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton,  
384 Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano,  
385 Jan Leike, and Ryan Lowe. Training language models to follow instructions with human  
386 feedback, 2022.
- 387 [22] Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. Squad: 100,000+ questions  
388 for machine comprehension of text. *arXiv preprint arXiv:1606.05250*, 2016.
- 389 [23] Rylan Schaeffer, Brando Miranda, and Sanmi Koyejo. Are emergent abilities of large language  
390 models a mirage? *arXiv preprint arXiv:2304.15004*, 2023.
- 391 [24] C Spearman. " general intelligence," objectively determined and measured. *The American*  
392 *Journal of Psychology*, 15(2):201–292, 1904.
- 393 [25] Aarohi Srivastava, Yinhan Deng, Nicholas Hay, Noam Shazeer, Ethan Paull, Doug Downey,  
394 Jonathan Duerig, Niranjana Sundaram, Andrew Bornstein, Harsh Trivedi, Kushal Doshi, Samyak  
395 Savarese, Nathaniel Daw, Jie Zhu, Marc Lanctot, Azalia Mirhoseini, Emilio Parisotto, Ruslan  
396 Salakhutdinov, Mohammad Shoeybi, Yuxuan Tian, Luke Hawkins-Hooker, William Fedus,  
397 Robyn Lingelbach, Deepak Pathak, Ilya Sutskever, and Igor Mordatch. Beyond the imitation  
398 game: Measuring and ensuring broad and robust capabilities in large language models. In  
399 *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics*, 2022.
- 400 [26] Lintang Sutawika et al. Eleutherai/lm-evaluation-harness: Major refactor, December 2023.  
401 URL <https://doi.org/10.5281/zenodo.10256836>.
- 402 [27] Bala Swaminathan, R Ekke Vaishali, and R subashriTS. Analysis of minimax algorithm using  
403 tic-tac-toe. 2020. URL <https://api.semanticscholar.org/CorpusID:228863323>.
- 404 [28] Paul Robinson Sweet. On language: The diversity of human language-structure and its influence  
405 on the mental development of mankind. by wilhelm von humboldt. translated by peter heath.  
406 *Historiographia Linguistica*, 16:387–392, 1989. URL <https://api.semanticscholar.org/CorpusID:170369059>.  
407
- 408 [29] Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy  
409 Liang, and Tatsunori B. Hashimoto. Stanford alpaca: An instruction-following llama model.  
410 [https://github.com/tatsu-lab/stanford\\_alpaca](https://github.com/tatsu-lab/stanford_alpaca), 2023.
- 411 [30] Oguzhan Topsakal and Jackson B. Harper. Benchmarking large language model (llm)  
412 performance for game playing via tic-tac-toe. *Electronics*, 2024. URL <https://api.semanticscholar.org/CorpusID:269225397>.  
413

- 414 [31] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timo-  
415 thée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez,  
416 Armand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation  
417 language models, 2023.
- 418 [32] A. M. TURING. I.—COMPUTING MACHINERY AND INTELLIGENCE. *Mind*, LIX  
419 (236):433–460, 10 1950. ISSN 0026-4423. doi: 10.1093/mind/LIX.236.433. URL <https://doi.org/10.1093/mind/LIX.236.433>.  
420
- 421 [33] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez,  
422 Lukasz Kaiser, and Illia Polosukhin. Attention is all you need, 2023.
- 423 [34] Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman.  
424 Glue: A multi-task benchmark and analysis platform for natural language understanding. *arXiv*  
425 *preprint arXiv:1804.07461*, 2018.
- 426 [35] Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill,  
427 Oyvind Levy, and Samuel R Bowman. Superglue: A stickier benchmark for general-purpose  
428 language understanding systems. *Advances in Neural Information Processing Systems*, 2019.
- 429 [36] D. Wechsler. *The measurement of adult intelligence*. Baltimore: Williams & Wilkins, 1939.
- 430 [37] Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can  
431 a machine really finish your sentence? In *Proceedings of the 57th Annual Meeting of the*  
432 *Association for Computational Linguistics*, 2019.
- 433 [38] Matej Zečević. Causal parrots: Large language models may talk causality but are not causal.  
434 *Unpublished Manuscript*, 2023.

435 **A Appendix / supplemental material**

436 **A.1 Move Mapping**

437 **A.1.1 Tic-Tac-Toe**

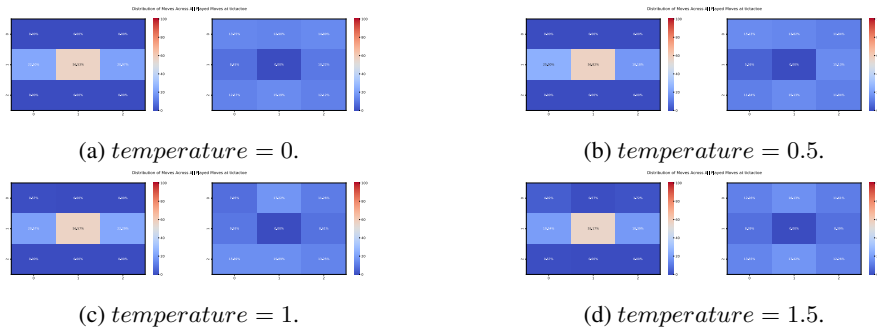


Figure 10: Heatmap of model GPT-3.5’s moves for the tic-tac-toe game.

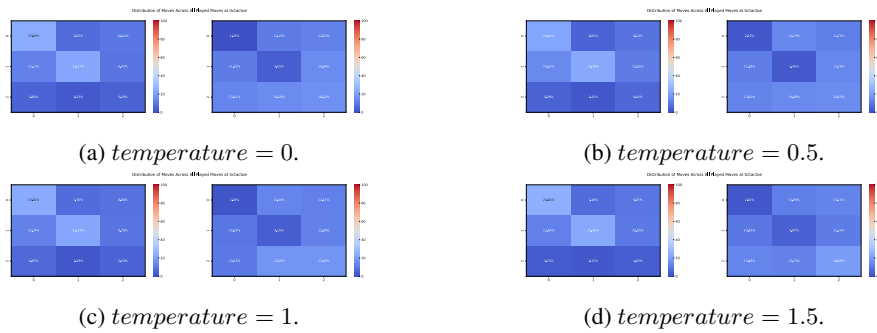


Figure 11: Heatmap of model GPT-4’s moves for the tic-tac-toe game.

438 **A.1.2 Connect-Four**

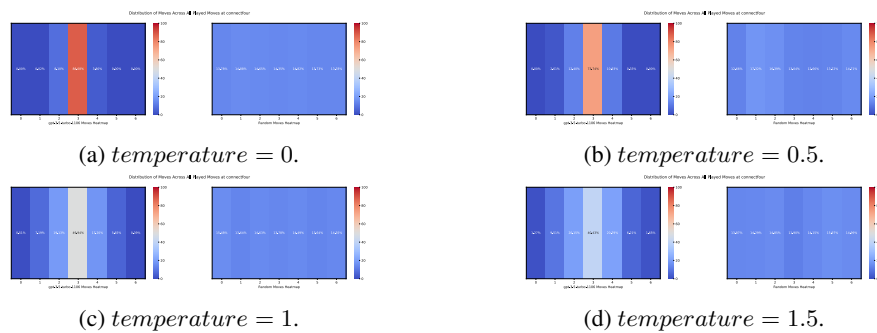


Figure 12: Heatmap of model GPT-3.5’s moves for the connect-four game.

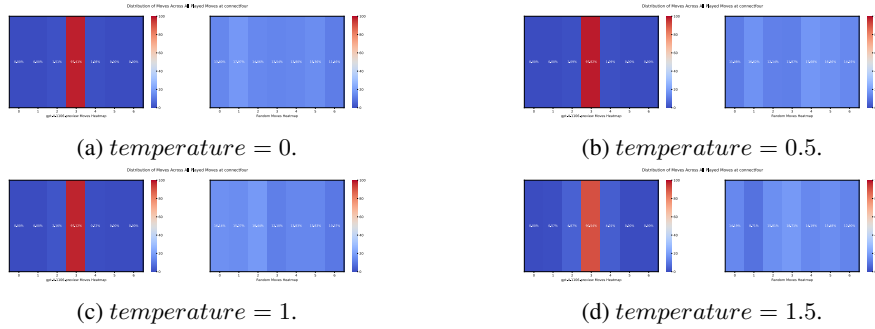


Figure 13: Heatmap of model GPT-4’s moves for the connect-four game.

439 **A.1.3 Battleship**

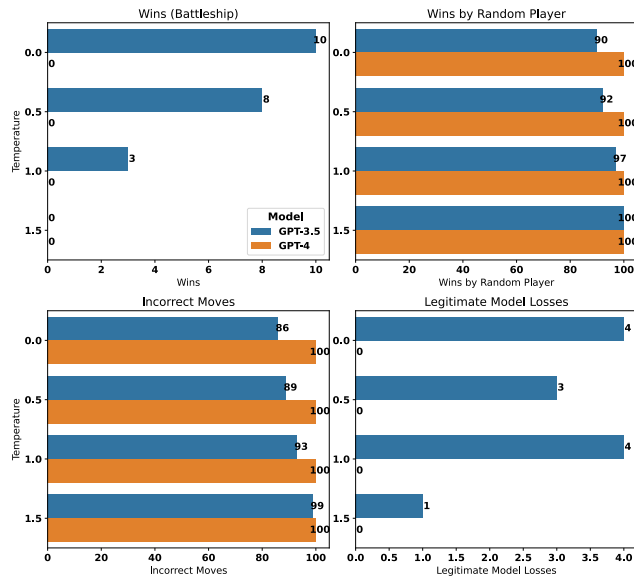


Figure 14: Battleship.

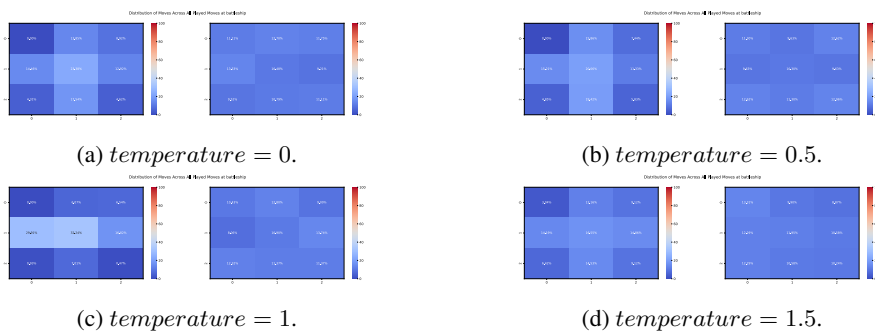


Figure 15: Heatmap of model GPT-3.5’s moves for the battleship game.

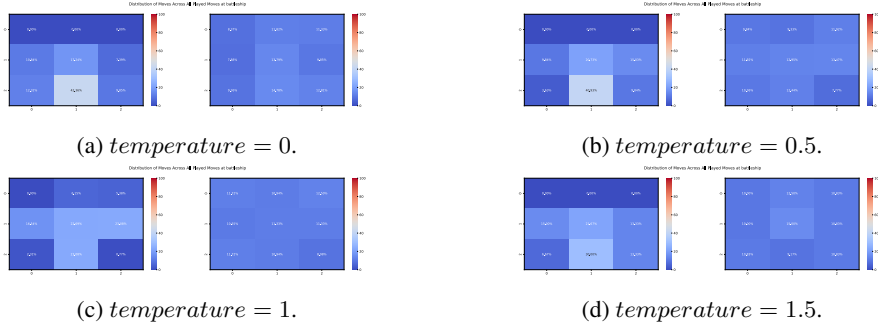


Figure 16: Heatmap of model GPT-4's moves for the battleship game.

440 **A.2 Shapes**

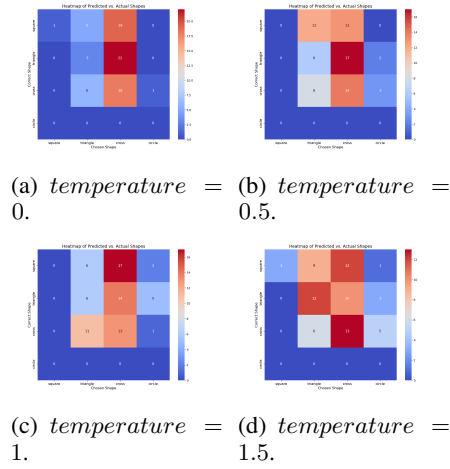


Figure 17: Heatmap of model GPT-3.5's moves for the shapes game.

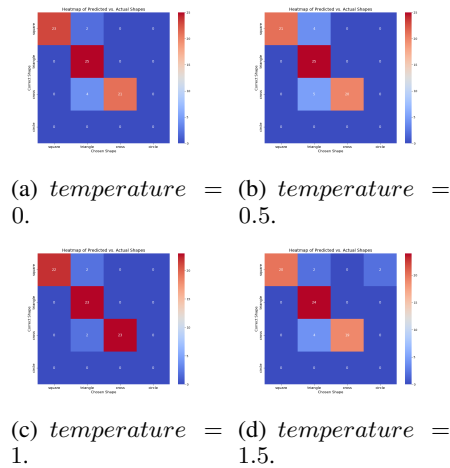


Figure 18: Heatmap of model GPT-4's moves for the shapes game.

441 **A.3 Prompting GPT About Optimal Play**

Game	Explanation
Tic-Tac-Toe	Tic-Tac-Toe is a two-player game played on a 3x3 grid. Each player takes turns marking a square with their symbol (X or O), aiming to get three of their symbols in a row, column, or diagonal. To play optimally, prioritize securing the center square and blocking opponent's winning moves.
Battleship	Battleship is a two-player game where players hide ships on a grid and take turns guessing their opponent's ship locations. The goal is to sink all of the opponent's ships. To play optimally, start by targeting areas with higher probabilities of containing a ship and strategically target adjacent squares after a hit to maximize efficiency.
Connect Four	Connect Four is a two-player game played on a 6x7 grid. Players drop colored discs into columns, aiming to connect four of their own discs in a row, column, or diagonal. To play optimally, prioritize creating your own winning formations while blocking opponent's potential winning moves.

Table 1: Optimal strategies for playing different games according to GPT-3.5.

Game	Explanation
Tic-Tac-Toe	Play your first X in a corner to maximize opportunities. If the opponent plays in the center, play the opposite corner. Block your opponent's potential winning moves and always look to create a line of three.
Battleship	Randomize ship placements and start by targeting the center of the grid. Use a checkerboard pattern for efficient searching. Once a ship is hit, focus on the surrounding squares to determine its orientation and sink it.
Connect Four	Start in the center column to maximize opportunities in all directions. Build threats vertically, horizontally, and diagonally, and block the opponent's forming lines. Create multiple threats to force the opponent into a defensive position.

Table 2: Optimal strategies for playing different games according to GPT-4.



## 442 **NeurIPS Paper Checklist**

### 443 **1. Claims**

444 Question: Do the main claims made in the abstract and introduction accurately reflect the  
445 paper's contributions and scope?

446 Answer: [\[Yes\]](#)

447 Justification: Yes, see sections 2 and 3, where we explore the delineated experiments and  
448 the ensuing results.

449 Guidelines:

- 450 • The answer NA means that the abstract and introduction do not include the claims  
451 made in the paper.
- 452 • The abstract and/or introduction should clearly state the claims made, including the  
453 contributions made in the paper and important assumptions and limitations. A No or  
454 NA answer to this question will not be perceived well by the reviewers.
- 455 • The claims made should match theoretical and experimental results, and reflect how  
456 much the results can be expected to generalize to other settings.
- 457 • It is fine to include aspirational goals as motivation as long as it is clear that these goals  
458 are not attained by the paper.

### 459 **2. Limitations**

460 Question: Does the paper discuss the limitations of the work performed by the authors?

461 Answer: [\[Yes\]](#)

462 Justification: See section 4, where we dive into some of the limitations of this work.

463 Guidelines:

- 464 • The answer NA means that the paper has no limitation while the answer No means that  
465 the paper has limitations, but those are not discussed in the paper.
- 466 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 467 • The paper should point out any strong assumptions and how robust the results are to  
468 violations of these assumptions (e.g., independence assumptions, noiseless settings,  
469 model well-specification, asymptotic approximations only holding locally). The authors  
470 should reflect on how these assumptions might be violated in practice and what the  
471 implications would be.
- 472 • The authors should reflect on the scope of the claims made, e.g., if the approach was  
473 only tested on a few datasets or with a few runs. In general, empirical results often  
474 depend on implicit assumptions, which should be articulated.
- 475 • The authors should reflect on the factors that influence the performance of the approach.  
476 For example, a facial recognition algorithm may perform poorly when image resolution  
477 is low or images are taken in low lighting. Or a speech-to-text system might not be  
478 used reliably to provide closed captions for online lectures because it fails to handle  
479 technical jargon.
- 480 • The authors should discuss the computational efficiency of the proposed algorithms  
481 and how they scale with dataset size.
- 482 • If applicable, the authors should discuss possible limitations of their approach to  
483 address problems of privacy and fairness.
- 484 • While the authors might fear that complete honesty about limitations might be used by  
485 reviewers as grounds for rejection, a worse outcome might be that reviewers discover  
486 limitations that aren't acknowledged in the paper. The authors should use their best  
487 judgment and recognize that individual actions in favor of transparency play an impor-  
488 tant role in developing norms that preserve the integrity of the community. Reviewers  
489 will be specifically instructed to not penalize honesty concerning limitations.

### 490 **3. Theory Assumptions and Proofs**

491 Question: For each theoretical result, does the paper provide the full set of assumptions and  
492 a complete (and correct) proof?

493 Answer: [\[NA\]](#)

494 Justification: We do not produce any theoretical results, rather we have made a benchmark  
495 and produce the experiments using said benchmark.

496 Guidelines:

- 497 • The answer NA means that the paper does not include theoretical results.
- 498 • All the theorems, formulas, and proofs in the paper should be numbered and cross-  
499 referenced.
- 500 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 501 • The proofs can either appear in the main paper or the supplemental material, but if  
502 they appear in the supplemental material, the authors are encouraged to provide a short  
503 proof sketch to provide intuition.
- 504 • Inversely, any informal proof provided in the core of the paper should be complemented  
505 by formal proofs provided in appendix or supplemental material.
- 506 • Theorems and Lemmas that the proof relies upon should be properly referenced.

#### 507 4. Experimental Result Reproducibility

508 Question: Does the paper fully disclose all the information needed to reproduce the main ex-  
509 perimental results of the paper to the extent that it affects the main claims and/or conclusions  
510 of the paper (regardless of whether the code and data are provided or not)?

511 Answer: [Yes]

512 Justification: See section 2.

513 Guidelines:

- 514 • The answer NA means that the paper does not include experiments.
- 515 • If the paper includes experiments, a No answer to this question will not be perceived  
516 well by the reviewers: Making the paper reproducible is important, regardless of  
517 whether the code and data are provided or not.
- 518 • If the contribution is a dataset and/or model, the authors should describe the steps taken  
519 to make their results reproducible or verifiable.
- 520 • Depending on the contribution, reproducibility can be accomplished in various ways.  
521 For example, if the contribution is a novel architecture, describing the architecture fully  
522 might suffice, or if the contribution is a specific model and empirical evaluation, it may  
523 be necessary to either make it possible for others to replicate the model with the same  
524 dataset, or provide access to the model. In general, releasing code and data is often  
525 one good way to accomplish this, but reproducibility can also be provided via detailed  
526 instructions for how to replicate the results, access to a hosted model (e.g., in the case  
527 of a large language model), releasing of a model checkpoint, or other means that are  
528 appropriate to the research performed.
- 529 • While NeurIPS does not require releasing code, the conference does require all submis-  
530 sions to provide some reasonable avenue for reproducibility, which may depend on the  
531 nature of the contribution. For example
  - 532 (a) If the contribution is primarily a new algorithm, the paper should make it clear how  
533 to reproduce that algorithm.
  - 534 (b) If the contribution is primarily a new model architecture, the paper should describe  
535 the architecture clearly and fully.
  - 536 (c) If the contribution is a new model (e.g., a large language model), then there should  
537 either be a way to access this model for reproducing the results or a way to reproduce  
538 the model (e.g., with an open-source dataset or instructions for how to construct  
539 the dataset).
  - 540 (d) We recognize that reproducibility may be tricky in some cases, in which case  
541 authors are welcome to describe the particular way they provide for reproducibility.  
542 In the case of closed-source models, it may be that access to the model is limited in  
543 some way (e.g., to registered users), but it should be possible for other researchers  
544 to have some path to reproducing or verifying the results.

#### 545 5. Open access to data and code

546 Question: Does the paper provide open access to the data and code, with sufficient instruc-  
547 tions to faithfully reproduce the main experimental results, as described in supplemental  
548 material?

549  
550  
551  
552  
553  
554  
555  
556  
557  
558  
559  
560  
561  
562  
563  
564  
565  
566  
567  
568  
569  
570  
571  
572  
573  
574  
575  
576  
577  
578  
579  
580  
581  
582  
583  
584  
585  
586  
587  
588  
589  
590  
591  
592  
593  
594  
595  
596  
597  
598  
599

Answer: [Yes]

Justification: We provide open access to our data and experiments through (GitHub Repository).

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

## 6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: We explicitly mention the temperature used in every plot and section 2.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

## 7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [No]

Justification: We have to rerun some of the experiments to recalculate these.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).

- 600 • It should be clear whether the error bar is the standard deviation or the standard error  
601 of the mean.
- 602 • It is OK to report 1-sigma error bars, but one should state it. The authors should  
603 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis  
604 of Normality of errors is not verified.
- 605 • For asymmetric distributions, the authors should be careful not to show in tables or  
606 figures symmetric error bars that would yield results that are out of range (e.g. negative  
607 error rates).
- 608 • If error bars are reported in tables or plots, The authors should explain in the text how  
609 they were calculated and reference the corresponding figures or tables in the text.

## 610 8. Experiments Compute Resources

611 Question: For each experiment, does the paper provide sufficient information on the com-  
612 puter resources (type of compute workers, memory, time of execution) needed to reproduce  
613 the experiments?

614 Answer: [Yes]

615 Justification: See section 2. We mention compute time, but all experiments are dependent  
616 on OpenAI's API.

617 Guidelines:

- 618 • The answer NA means that the paper does not include experiments.
- 619 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,  
620 or cloud provider, including relevant memory and storage.
- 621 • The paper should provide the amount of compute required for each of the individual  
622 experimental runs as well as estimate the total compute.
- 623 • The paper should disclose whether the full research project required more compute  
624 than the experiments reported in the paper (e.g., preliminary or failed experiments that  
625 didn't make it into the paper).

## 626 9. Code Of Ethics

627 Question: Does the research conducted in the paper conform, in every respect, with the  
628 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

629 Answer: [Yes]

630 Justification:

631 Guidelines:

- 632 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 633 • If the authors answer No, they should explain the special circumstances that require a  
634 deviation from the Code of Ethics.
- 635 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-  
636 eration due to laws or regulations in their jurisdiction).

## 637 10. Broader Impacts

638 Question: Does the paper discuss both potential positive societal impacts and negative  
639 societal impacts of the work performed?

640 Answer:[Yes]

641 Justification: See sections 4 and 5 where we go over the implications of our results in the  
642 context of LLM interpretation.

643 Guidelines:

- 644 • The answer NA means that there is no societal impact of the work performed.
- 645 • If the authors answer NA or No, they should explain why their work has no societal  
646 impact or why the paper does not address societal impact.
- 647 • Examples of negative societal impacts include potential malicious or unintended uses  
648 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations  
649 (e.g., deployment of technologies that could make decisions that unfairly impact specific  
650 groups), privacy considerations, and security considerations.

- 651 • The conference expects that many papers will be foundational research and not tied  
652 to particular applications, let alone deployments. However, if there is a direct path to  
653 any negative applications, the authors should point it out. For example, it is legitimate  
654 to point out that an improvement in the quality of generative models could be used to  
655 generate deepfakes for disinformation. On the other hand, it is not needed to point out  
656 that a generic algorithm for optimizing neural networks could enable people to train  
657 models that generate Deepfakes faster.
- 658 • The authors should consider possible harms that could arise when the technology is  
659 being used as intended and functioning correctly, harms that could arise when the  
660 technology is being used as intended but gives incorrect results, and harms following  
661 from (intentional or unintentional) misuse of the technology.
- 662 • If there are negative societal impacts, the authors could also discuss possible mitigation  
663 strategies (e.g., gated release of models, providing defenses in addition to attacks,  
664 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from  
665 feedback over time, improving the efficiency and accessibility of ML).

## 666 11. Safeguards

667 Question: Does the paper describe safeguards that have been put in place for responsible  
668 release of data or models that have a high risk for misuse (e.g., pretrained language models,  
669 image generators, or scraped datasets)?

670 Answer: [NA]

671 Justification: There is no production of models or data that would pose risk.

672 Guidelines:

- 673 • The answer NA means that the paper poses no such risks.
- 674 • Released models that have a high risk for misuse or dual-use should be released with  
675 necessary safeguards to allow for controlled use of the model, for example by requiring  
676 that users adhere to usage guidelines or restrictions to access the model or implementing  
677 safety filters.
- 678 • Datasets that have been scraped from the Internet could pose safety risks. The authors  
679 should describe how they avoided releasing unsafe images.
- 680 • We recognize that providing effective safeguards is challenging, and many papers do  
681 not require this, but we encourage authors to take this into account and make a best  
682 faith effort.

## 683 12. Licenses for existing assets

684 Question: Are the creators or original owners of assets (e.g., code, data, models), used in  
685 the paper, properly credited and are the license and terms of use explicitly mentioned and  
686 properly respected?

687 Answer: [Yes]

688 Justification: Authors are credited and the license is made available in GitHub Repository.

689 Guidelines:

- 690 • The answer NA means that the paper does not use existing assets.
- 691 • The authors should cite the original paper that produced the code package or dataset.
- 692 • The authors should state which version of the asset is used and, if possible, include a  
693 URL.
- 694 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 695 • For scraped data from a particular source (e.g., website), the copyright and terms of  
696 service of that source should be provided.
- 697 • If assets are released, the license, copyright information, and terms of use in the  
698 package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets)  
699 has curated licenses for some datasets. Their licensing guide can help determine the  
700 license of a dataset.
- 701 • For existing datasets that are re-packaged, both the original license and the license of  
702 the derived asset (if it has changed) should be provided.

703 • If this information is not available online, the authors are encouraged to reach out to  
704 the asset’s creators.

705 **13. New Assets**

706 Question: Are new assets introduced in the paper well documented and is the documentation  
707 provided alongside the assets?

708 Answer: [Yes]

709 Justification: See section 2.

710 Guidelines:

- 711 • The answer NA means that the paper does not release new assets.
- 712 • Researchers should communicate the details of the dataset/code/model as part of their  
713 submissions via structured templates. This includes details about training, license,  
714 limitations, etc.
- 715 • The paper should discuss whether and how consent was obtained from people whose  
716 asset is used.
- 717 • At submission time, remember to anonymize your assets (if applicable). You can either  
718 create an anonymized URL or include an anonymized zip file.

719 **14. Crowdsourcing and Research with Human Subjects**

720 Question: For crowdsourcing experiments and research with human subjects, does the paper  
721 include the full text of instructions given to participants and screenshots, if applicable, as  
722 well as details about compensation (if any)?

723 Answer: [NA]

724 Justification: We do not involve people directly in our experiments.

725 Guidelines:

- 726 • The answer NA means that the paper does not involve crowdsourcing nor research with  
727 human subjects.
- 728 • Including this information in the supplemental material is fine, but if the main contribu-  
729 tion of the paper involves human subjects, then as much detail as possible should be  
730 included in the main paper.
- 731 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,  
732 or other labor should be paid at least the minimum wage in the country of the data  
733 collector.

734 **15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human  
735 Subjects**

736 Question: Does the paper describe potential risks incurred by study participants, whether  
737 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)  
738 approvals (or an equivalent approval/review based on the requirements of your country or  
739 institution) were obtained?

740 Answer: [NA]

741 Justification: The paper does not involve crowdsourcing nor research with human subjects.

742 Guidelines:

- 743 • The answer NA means that the paper does not involve crowdsourcing nor research with  
744 human subjects.
- 745 • Depending on the country in which research is conducted, IRB approval (or equivalent)  
746 may be required for any human subjects research. If you obtained IRB approval, you  
747 should clearly state this in the paper.
- 748 • We recognize that the procedures for this may vary significantly between institutions  
749 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the  
750 guidelines for their institution.
- 751 • For initial submissions, do not include any information that would break anonymity (if  
752 applicable), such as the institution conducting the review.