

# INVESTIGATING THE LINK BETWEEN REPRESENTATIONAL SIMILARITY AND MODEL INTERACTIONS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Researchers have shown that neural similarity among humans predicts social closeness and cooperative success, whereas innovation often emerges from interactions among dissimilar individuals. We investigate whether these principles extend to artificial intelligence by examining interactions between large language models. In our experiments, 276 model pairs interact across eight games spanning both cooperation and novelty. We find that pairs with more similar representation spaces achieve significantly higher cooperation but exhibit reduced novelty and creativity. [The effects of representational similarity on cooperation and novelty remain robust even after isolating other factors such as performance disparity and model size.](#) We also find that similarity in the early layers consistently exhibits the strongest effect across games, compared to the middle and later layers. [This suggests that a central factor underlying the observed trend is the extent to which the two models share lexical and semantic grounding.](#) These findings suggest that representational similarity can be an important consideration in multi-agent system design.

## 1 INTRODUCTION

The deployment of multiple large language models (LLMs) in multi-turn, multi-agent interactions has progressed rapidly from concept to practice, with recent investigations in applications to social simulations (Park et al., 2023; Xie et al., 2024; Zhou et al., 2023), coding (Wu et al., 2024a; Ishibashi & Nishimura, 2024), and a range of creative tasks such as writing, brainstorming, and scientific idea generation (Chen et al., 2023; Fukumura & Ito, 2025; Su et al., 2025). In many collaborative tasks, prior work has found that interaction between multiple agents facilitates stronger performance than single-agent systems (Talebirad & Nadiri, 2023; Zhuge et al., 2023). Beyond treating multi-agent systems as tools, some have even proposed evolving LLMs through multi-agent interaction (Lai et al., 2024; Eisenstein et al., 2025; Wu et al., 2025).

On the other hand, by their very nature, multi-agent systems are more complex than single-agent systems, increasing the potential for unexpected behaviors (Piatti et al., 2024; Hammond et al., 2025; de Witt, 2025). One central concern is whether agents can reliably cooperate with one another, which is critical for the safe and reliable deployment of multi-agent systems (Piedrahita et al., 2025). Being able to understand and predict the dynamics of multi-agent systems is therefore essential. Yet, most efforts to date have focused on single-agent cases, while studies of multi-agent systems have primarily focused on output-level behaviors rather than internal mechanisms.

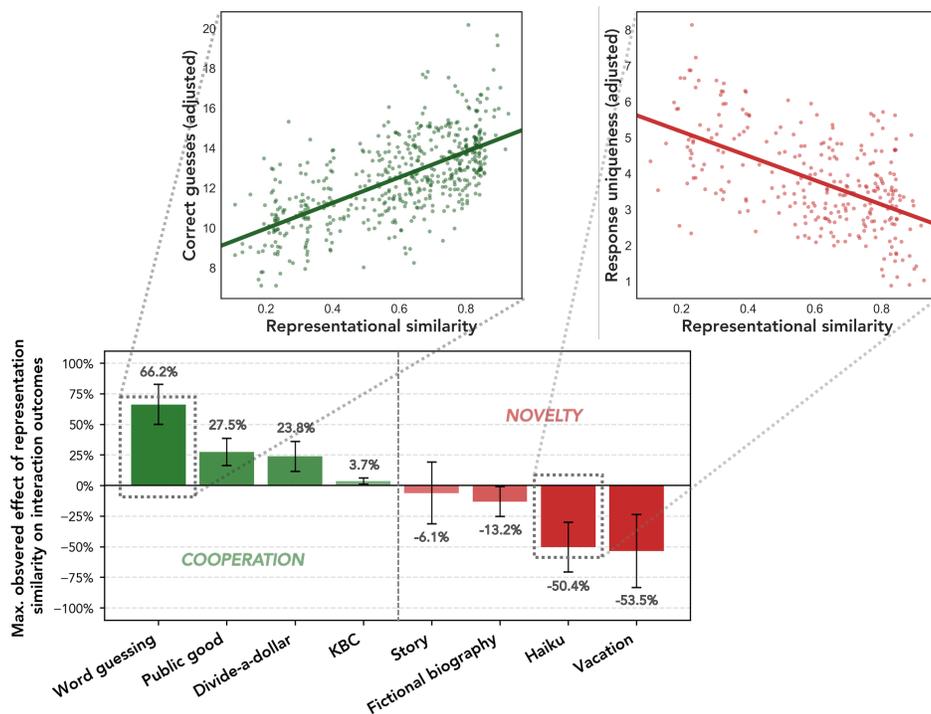
This work provides an initial exploration of multi-agent interaction through the lens of representational alignment. Specifically, we ask:

*What is the relationship between representational similarity and interactive behavior of models?*

Evidence from neuroscience and social sciences suggests that similar neural responses among humans are significantly associated with their social closeness and cooperative performance (Parkinson et al., 2018; Thornton & Mitchell, 2017; Shen et al., 2025b; Reiner et al., 2021), while interaction between dissimilar individuals often sparks innovation (Hewlett et al., 2013; Østergaard et al., 2011). Analogously, we hypothesize that *models with higher representational similarity are more likely to cooperate and predict one another, but exhibit reduced collective novelty and creativity.*

To test this, we conduct experiments involving 276 model pairs spanning 23 open-weight LLMs from eight model families. Specifically, we examine cooperation through four games: word guessing, public good,

054 divide-a-dollar, and the Keynesian Beauty Contest (KBC); and assess creativity and novelty through four  
 055 generative tasks: story writing, fictional biography, haiku composition, and vacation benefit brainstorming.  
 056



080 Figure 1: The effect of representational similarity on each game outcome. Representational similarity is  
 081 quantified using linear Centered Kernel Alignment (CKA) (Kornblith et al., 2019) with WikiText (Merity  
 082 et al., 2016). In the bar graph, the effect size reflects the relative change (%) in predicted outcomes  
 083 between the lowest and the highest observed value of representational similarity, with error bars indicating  
 084 95% confidence intervals. In the scatter plots, each point represents a model pair, and the  $y$  values  
 085 are adjusted via mixed-effects regression to control for model-specific tendencies, thereby isolating the effect  
 086 of representational similarity on interaction outcomes. Overall, greater similarity corresponds to higher  
 087 cooperation but lower novelty.

088 Our experiments reveal that representational similarity is a strong predictor of interactive outcomes.  
 089 Figure 1 illustrates how these outcomes vary with increasing internal similarity across scenarios:  
 090 cooperation performance rises significantly as representational similarity increases. For example, in the  
 091 word-guessing game where one player attempts to identify their partner’s secret word, correct guesses  
 092 increase by roughly 66.2% (relative change) as representational similarity rises from the minimum to  
 093 the maximum observed values. By contrast, novelty declines consistently across the four creative tasks,  
 094 though the magnitude and statistical significance vary. **Even when accounting for other factors such as**  
 095 **performance disparity, representational similarity shows a strong effect on cooperation and novelty.** These  
 096 findings suggest a likely tradeoff: model pairs with higher representational similarity tend to cooperate  
 097 better, but also manifest reduced collective novelty. These results provide new insights into the design  
 098 of multi-agent systems, where single-model deployment is currently the dominant paradigm.

## 100 2 RELATED WORK

101  
 102 **Multi-Agent Systems and Their Behaviors.** A growing body of literature examines how LLMs interact  
 103 and cooperate within multi-agent systems (Piatti et al., 2024; Lai et al., 2024; Li & Shirado, 2025; Piedrahita  
 104 et al., 2025; Wu et al., 2024b; Zhu et al., 2025; Kim, 2025). To study cooperative behaviors, these works  
 105 commonly employ economic games such as the public good game (Hauert et al., 2006) and the tragedy of the  
 106 commons (Hardin, 1968), both of which have long served as canonical paradigms for analyzing cooperation.  
 107 Lai et al. (2024) found that agents become substantially more cooperative in the public good game, after en-  
 gaging in multi-round, free-form interactions with peers. Wu et al. (2024b) further showed that LLM agents

108 can develop emergent cooperative strategies across various environments, even in settings that are compet-  
109 itive or only partially cooperative. In contrast, several studies focusing on reasoning models have reported  
110 reduced cooperation, where these models often act as free riders in economic games (Li & Shirado, 2025;  
111 Piedrahita et al., 2025). More relevant to our study, Kim (2025) found that agents exhibit self-awareness,  
112 such that the Nash equilibrium in a game converges faster when they play the game with themselves.

113 Another line of research compares the performance of multi-agent systems with that of single-agent  
114 systems (Chen et al., 2023; Lai et al., 2024; Li et al., 2024; Du et al., 2023; Su et al., 2025). For example,  
115 AutoAgents (Chen et al., 2023) demonstrated that sharing diverse perspectives among agents can yield  
116 more creative novel writing. Lai et al. (2024) similarly found that agents that undergo free-form interaction  
117 exhibit enhanced creativity in sentence-generation tasks. Building on these prior studies, we investigate  
118 cooperative behaviors and creativity of multi-agent systems.

119 **Neural Similarity as a Predictor of Interaction in Humans and Models.** In neuroscience, similar  
120 neural responses between humans significantly predict social dynamics: Parkinson et al. (2018) found  
121 that friendship formation is predicted by similarity of neural response patterns to videos, as measured  
122 by fMRI. Shen et al. (2025b) extended this analysis to the similarity of neural activations during  
123 story-listening. Others have also shown that consistent neural activity patterns appear among personally  
124 familiar individuals (Thornton & Mitchell, 2017; Hyon et al., 2020). Such consistent findings suggest  
125 that greater neural similarity may facilitate stronger social bonds. Moreover, a related body of research  
126 has investigated the relationship between neural similarity and cooperative performance (Cui et al., 2012;  
127 Hu et al., 2017; Reiner et al., 2021; Réveillé et al., 2024), where it has been consistently found that higher  
128 interbrain synchrony is positively associated with cooperation.

129 As AI models scale in size and improve in performance, their internal representations increasingly align  
130 with human neural activity patterns (Goldstein et al., 2022; Schrimpf et al., 2021; Caucheteux & King,  
131 2022; Shen et al., 2025a; Gurnee et al., 2023; Zhou et al., 2025). For example, Caucheteux & King  
132 (2022) showed that language algorithms predicting words exhibit representational patterns similar to brain  
133 responses to sentences. Shen et al. (2025a) reported a strong correlation between brain-model similarity  
134 scores and model performance across both language models and vision models. This growing evidence  
135 of brain-model alignment motivates our central hypothesis. It raises the possibility that principles observed  
136 in human cognition and social behavior may also extend to advanced artificial systems.

137 **Representational Similarity in Neural Networks.** Researchers have long sought to understand the  
138 behavior of neural networks by comparing their internal representations. A variety of metrics for  
139 representational similarity in artificial neural networks have been proposed (Kornblith et al., 2019;  
140 Hotelling, 1992; Morcos et al., 2018; Raghu et al., 2017; Kriegeskorte et al., 2008). One widely used  
141 method is Centered Kernel Alignment (CKA; Kornblith et al., 2019), which enables comparison of  
142 representations between models regardless of their architecture or layer count. Several studies have  
143 investigated the nuances of applying these metrics. Ding et al. (2021) evaluated the sensitivity of similarity  
144 measures to changes in model behavior and showed that different metrics exhibit distinct weaknesses.  
145 Moschella et al. (2022) demonstrated that representational similarity can serve as a strong predictor of  
146 model performance, in tasks such as classification with vision models.

147 Beyond standard similarity metrics, new approaches have been proposed for comparing representation  
148 spaces. For example, *model stitching*—connecting two neural networks—has been argued to capture aspects  
149 of representational structure that metrics like CKA can not (Lenc & Vedaldi, 2015; Bansal et al., 2021). In  
150 this view, models with greater similarity are expected to achieve higher stitching success. **However, model  
151 stitching is substantially more expensive and less scalable than other representational similarity metrics, as  
152 it requires training a connector layer between two models.** Hacoen et al. (2020) proposed comparing the  
153 similarity of classification predictions in vision models as an alternative perspective on model comparison.

154 **Diversity, Creativity, and Collective Intelligence.** Behavioral research on innovation finds that higher  
155 diversity within a group of collaborators leads to increased novelty in their creations. For example, Uzzi  
156 et al. (2013) analyzed millions of scientific papers and found that the highest-impact science often arises  
157 from groups that combined existing research in novel ways. Page et al. (2019) formalizes this and proves  
158 that functionally diverse groups outperform homogeneous ones on complex problems, demonstrating  
159 superior problem-solving, innovation, and prediction accuracy. Similarly, Paulus (2000) showed that the  
160 effectiveness of brainstorming depends on cognitive diversity—that is, differences in how individuals  
161 perceive and think. Our results empirically explore this human-inspired principle: Does representational  
diversity within sets of LLM agents predict greater novelty in multi-agent creative tasks?

### 3 REPRESENTATIONAL SIMILARITY OF LLMs

To test our hypothesis—whether there is a relationship between representational similarity and interactive behavior of models—we first need a way to measure representational similarity of LLMs. In this section, we describe how we compute this similarity. *It is important to note that CKA computation is conducted independently of model interaction.*

#### 3.1 SIMILARITY METRICS

Representational similarity quantifies how similarly two neural models embed the same inputs. Measuring similarity involves two steps: **1)** extracting representational vectors from each model using a probe dataset (i.e., a set of prompts) and **2)** computing a similarity score between the extracted representations using a metric.

**Step 1. Extracting representations.** The first step can be formalized as follows. The probe dataset  $\mathcal{D} \subset \mathcal{X}$  contains  $m$  texts  $x \in \mathcal{X}$ , where  $\mathcal{X}$  is the set of all possible texts. Thus,  $\mathcal{D} = \{x_i\}_{i=1}^m$ . For a neural model with parameters  $\theta$ , we define  $f_\theta^k: \mathcal{X} \rightarrow \mathbb{R}^n$  as the mapping from a text  $x \in \mathcal{X}$  to an  $n$ -dimensional activation at the  $k$ -th layer, where  $1 \leq k \leq l$  and the model has  $l$  layers. Stacking the embeddings for all  $x \in \mathcal{D}$  yields a matrix  $R_\theta^k \in \mathbb{R}^{m \times n}$ , with the  $i$ -th row equal to  $f_\theta^k(x_i)$ .

**Step 2. Computing similarity.** The next step is to compute similarity between the representational spaces  $\{R_{\theta_1}^i\}_{1 \leq i \leq l_1}$  and  $\{R_{\theta_2}^j\}_{1 \leq j \leq l_2}$ , for two models with parameters  $\theta_1, \theta_2$ , depths  $l_1, l_2$ , and hidden dimensions  $n_1, n_2$ . A variety of similarity metrics ( $M$ ) have been proposed, including Centered Kernel Alignment (CKA; Kornblith et al., 2019), Canonical Correlation Analysis (CCA; Hotelling, 1992; Morcos et al., 2018), Singular Vector Canonical Correlation Analysis (SVCCA; Raghu et al., 2017), and Representational Similarity Analysis (RSA; Kriegeskorte et al., 2008). Each defines a function  $M: \mathbb{R}^{m \times n_1} \times \mathbb{R}^{m \times n_2} \rightarrow \mathbb{R}$  that takes two matrices (i.e.,  $R_{\theta_1}^i$  and  $R_{\theta_2}^j$ ) as input.

We use CKA (Kornblith et al., 2019) given its popularity of use in prior work (Ciernik et al., 2024; Shen et al., 2025a; Liu et al., 2025). CKA enables the comparison between two models with different architectures and different numbers of layers. There are four common CKA variants: linear CKA, RBF CKA, unbiased linear CKA, and unbiased RBF CKA. These CKA values range in  $[0, 1]$ , with higher values indicating greater similarity. Following prior work (Liu et al., 2025; Shen et al., 2025a; Zou et al., 2023; Raffel et al., 2020), we obtain  $f_{\theta_1}^i(x)$  and  $f_{\theta_2}^j(x)$  from the activation of the last token of each input  $x$  at their respective layers. CKA scores are then calculated for every layer pair of the two models. That is,  $\text{CKA}(R_{\theta_1}^i, R_{\theta_2}^j)$  for all  $1 \leq i \leq l_1$  and  $1 \leq j \leq l_2$ , producing an  $l_1 \times l_2$  grid of scores.

To summarize similarity with a single score per model pair, there are multiple approaches. The first approach is to average the CKA scores (i.e., global average):  $\frac{\sum_{i,j} \text{CKA}(R_{\theta_1}^i, R_{\theta_2}^j)}{l_1 \cdot l_2}$ . This captures overall similarity between all layers of the two models. Please note that identical model pairs can score below 1, since off-diagonal layer pairs ( $i \neq j$ ) yield values less than 1. An alternative summary measure of CKA is the layer-wise maximum-aligned average, which captures how well each layer aligns with its best-matching layer in the other model. That is,

$$\frac{1}{2} \times \left( \frac{\sum_i \max_j \text{CKA}(R_{\theta_1}^i, R_{\theta_2}^j)}{l_1} + \frac{\sum_j \max_i \text{CKA}(R_{\theta_1}^i, R_{\theta_2}^j)}{l_2} \right).$$

With this measure, identical model pairs always achieve 1, since each layer’s best match is itself and  $\text{CKA}(R_\theta^i, R_\theta^i) = 1$ .

We observe consistent trends between representational similarity and interactive behavior across both aggregation methods and all four CKA variants. Unless otherwise noted, CKA refers to the global averages of linear CKA. Results for other variants appear in Appendix E.

#### 3.2 PROBE DATASET

A recent study (Ciernik et al., 2024) shows that representational similarity can depend on the choice of probe dataset. To examine whether the relationship between representational similarity and model interactions depends on probe dataset, we use four probe datasets spanning different domains, from which we compute a CKA score for each: WikiText (Merity et al., 2016) for general language, GSM8K (Cobbe

et al., 2021) and MATH (Hendrycks et al., 2021) for mathematics, and TruthfulQA (Lin et al., 2021) for truthfulness. From WikiText and MATH, we randomly sample 1000 prompts each. For GSM8K, we use the entire test set (1319 prompts), and for TruthfulQA, we use the full dataset (817 prompts).

### 3.3 REPRESENTATIONAL SIMILARITY RANGE

We consider 23 open-weight LLMs spanning eight families and sizes ranging from 1B to 72B parameters, yielding 276 model pairs. The full model list is provided in Table 1 in Appendix. These models exhibit a wide range of representational similarity. For example, using the global average score with WikiText as the probe dataset, values range from 0.106 (gemma-3-4b-it vs. gemma-3-12b-it) to 0.92 (phi-4 vs. phi-4). Using the average of maximum-aligned scores, values range from 0.288 (gemma-3-4b-it vs. gemma-3-12b-it) to 1 (for all identical model pairs). The complete set of CKA scores for all 276 pairs is shown in Figures 6 and 7. We find that the Gemma family (Kamath et al., 2025) generally exhibits lower similarity to other models, while pairs within the same family tend to show higher similarity. Figure 8 also reports correlations across different CKA variants, where similarities computed with GSM8K and WikiText display relatively lower agreement.

## 4 COOPERATION INCREASES WHEN SIMILAR MODELS MEET

Here, we define cooperation broadly to refer to behaviors in which agents align their actions or expectations to achieve mutually beneficial outcomes (Bowles & Gintis, 2003; Grice, 1975). Building on evidence that greater interbrain synchrony among humans is linked to increased cooperation, we test the hypothesis that model pairs with higher representational similarity will demonstrate increased cooperative behavior.

### 4.1 GAME SETTINGS & ANALYSIS

We use four game settings that involve cooperation: word guessing (Clark & Wilkes-Gibbs, 1986), public goods (Hauert et al., 2006), divide-a-dollar (Kalai, 1977), and the Keynesian Beauty Contest (KBC) (Duffy & Nagel, 1997). The word-guessing game is a form of referential communication, in which a speaker provides a clue referring to a target word and a listener attempts to identify the target based on that clue. Referential communication has long been used to examine how players coordinate and collaborate (Clark & Wilkes-Gibbs, 1986). The latter three games have been widely adopted in economics and social science to study cooperative and coordination dynamics among humans. Moreover, these four games have commonly been employed to study the behaviors of AI models (Tang et al., 2024; Lai et al., 2024; Wu et al., 2024b; Piedrahita et al., 2025; Li & Shirado, 2025; Huang et al., 2024; Kim, 2025). Based on the prior literature, we adopt these four games in our study. The following presents a description of each game rule along with associated outcome metrics, which capture the extent to which the two agents cooperated with one another during a game. Please refer to Appendix C for the prompts used in each game, along with examples.

- Word Guessing:** In the game, one player chooses their own target word that begins with a given letter (“a” to “z”) and provides a one-word hint to the other player. Here, the player is instructed to make the hint different from the target word. The other player should guess that secret word based on the hint and the given starting letter. Each round is one-shot and independent. We use the number of correct guesses over 26 rounds, one for each letter in the alphabet, as the outcome metric. Please note that it is not trivial to expect more similar models to perform better in this setting, because the game itself is asymmetric. The roles of the clue giver and the guesser are fundamentally distinct, and prior work in linguistics and cognitive science has shown that producing an informative cue and interpreting one are not mirror tasks (Hendriks & Koster, 2010; Hendriks, 2014; Mayol, 2018; Hendriks, 2016). This implies that even for a single individual, generating a hint and inferring the correct target from a self-produced hint rely on different underlying mechanisms. Consequently, success in the word-guessing game reflects cooperative behavior that goes beyond simple alignment or representational homogeneity.
- Public Good:** The game repeats for five rounds. At the beginning of the game, each agent begins with \$100 of their own money and decides how much to contribute to a public pot every round. After their contribution is collected into the public pot each round, the value increases by 30% and is evenly redistributed back to each agent. We use each agent’s total asset value accumulated over five rounds as the outcome metric. In this game, a non-cooperative strategy is to free-ride by contributing nothing, which maximizes an individual’s own total assets but substantially harms the collective welfare. Thus, the game captures how agents make a decision between cooperative and non-cooperative behavior across the five rounds.

- **Divide a Dollar:** The game repeats for five rounds. Each round, \$1 is available, and players should demand how much of the \$1 they want. If the total amount requested is not above \$1, players receive the amount they requested. If the total amount requested exceeds \$1, agents don't receive anything. We use each agent's total asset value accumulated over five rounds as the outcome metric. [In this game, aggressively demanding larger amounts without considering the other player's request reduces the likelihood that either player receives anything. Thus, the game captures how well agents engage in cooperative decision-making over five rounds while taking the other player's choices into account.](#)
- **KBC:** The game repeats for five rounds. At the beginning of each round, players choose a number between 0 and 100, guessing the closest number to  $2/3$  of the average of the numbers from both agents. The score is based on how close their number is to  $2/3$  of the average:  $100 - |\text{their number} - 2/3 \times \text{average}|$ . We use the total score of each player over the five rounds as the outcome metric. [A higher score reflects a player's ability to anticipate the other player's choice. Thus, the game captures how effectively agents engage in recursive reasoning about the other player's reasoning process and selected numbers.](#)

In all games except the word guessing game, where each round is one-shot and independent, players are shown the other's choice and reasoning at the end of each round. A higher game outcome value indicates stronger cooperation in that game. For example, in the word guessing game, performance depends on how accurately each agent guesses the other's secret words—reflecting their ability to interpret their partner and infer unknown information. In the public goods game, achieving high returns requires both cooperation and alignment: purely selfish strategies yield low payoffs, and exploitation due to misunderstanding also reduces outcomes.

We evaluate all 276 possible pairs of the 23 models listed in Table 1. Because the word guessing game is asymmetric, we consider ordered pairs, resulting in 529 model pairings. Each pair interacts across all four games, with temperature set to 0.7<sup>1</sup> and at least 4 independent samples collected per pair for each game. The average game outcome for each model is presented in Figure 9 in Appendix D.

To analyze the relationship between representational similarity and interaction outcomes, we fit a mixed-effects linear regression model (Bates et al., 2015). In our experimental setup, using a simple linear regression or Pearson correlation would be inappropriate because these tests assume independent data points, whereas our setup produces multiple samples per model pair, and each model appears in multiple pairs. Mixed-effects regression is the standard approach for handling such non-independence (Brown, 2021). In particular, it allows us to account for variance attributable to individual models (e.g., differences in capability) by including model-specific random effects, thereby isolating the effect of representational similarity on interactive outcomes. Specifically, we estimate the following mixed-effects regression:

$$Y_{ij} = \alpha + \beta \cdot \text{CKA}_{ij} + u_i + v_j + \epsilon_{ij},$$

where  $Y_{ij}$  is the interactive outcome of interest, and  $\text{CKA}_{ij}$  is the similarity measure between models  $i$  and  $j$ . The terms  $u_i$  and  $v_j$  represent random effects associated with models  $i$  and  $j$ , respectively, where these terms capture unobserved heterogeneity at the level of model  $i$  and model  $j$ , respectively. Lastly,  $\epsilon_{ij} \sim N(0, \sigma_\epsilon^2)$  is an error term.

*To evaluate whether similarity predicts the interactive outcome, the key quantities are the estimated slope of  $\text{CKA}_{ij}$  (i.e.,  $\beta$ ) and its statistical significance. We therefore report  $\beta$  with its  $p$ -value throughout the paper.*

## 4.2 DOES REPRESENTATIONAL SIMILARITY PREDICT COOPERATION?

Our results reveal that representational similarity can predict cooperative outcomes. In the word guessing game, correct guesses increase by approximately 88.2, 58.0, 63.3, and 59.4% (relative changes) for each unit increase in representational similarity (i.e., from 0 to 1) measured with WikiText, GSM8K, MATH, and TruthfulQA, respectively. In the public good game, each player's total asset value rises by 34.8, 32.4, 33.0, and 29.8% across the four probe datasets. For divide-a-dollar, asset values increase by 29.9, 15.3, 16.2, and 13.7%. Finally, in KBC, scores increase significantly but modestly—4.5, 4.7, 4.2, and 3.4%. All effects are found to be statistically significant (Figure 2). Among the four games, KBC shows the weakest effect. This is expected: the game has a unique Nash equilibrium in which both players always choose zero, which makes the optimal strategy fixed regardless of representational similarity. For instance, we observe that a certain model such as GPT-OSS-20B always chooses 0 regardless of the partner's decision. Nevertheless, even here, we observe a significant upward trend with increasing similarity.

<sup>1</sup>For robustness, we rerun the experiments with a temperature of 0.3 and observe that the results remain consistent (see Appendix F).

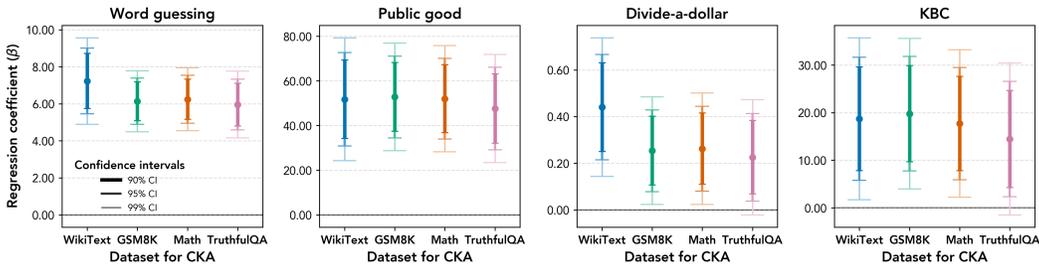


Figure 2: Regression coefficient of representational similarity on game outcomes. Error bars denote 90%, 95%, and 99% confidence intervals. Across games and datasets, the graphs show a positive effect of similarity on outcomes, with WikiText-based similarity exhibiting the strongest effect.

The pattern persists across probe datasets, implying its generalizability. Moreover, we find no difference in effect size across datasets (please refer to Figure 2). This contrasts with a previous finding Ciernik et al. (2024), which showed that the correspondence between representational similarity and task behavior depends on the dataset. The same trend holds across other CKA variants as well (see Appendix E.1).

### 5 NOVELTY DECREASES WHEN SIMILAR MODELS MEET

Next, we examine whether representational similarity predicts novelty in collaborative generative tasks. For this purpose, we adapt four tasks—story writing, fictional biography, haiku composition, and vacation benefit brainstorming—from NoveltyBench (Zhang et al., 2025), a benchmark originally designed to evaluate an individual model’s ability to produce high-quality and original ideas. Because NoveltyBench tasks are defined for single-agent settings, we extend them to the multi-agent case: each of the two models first generates a set of brainstorming ideas, after which each model produces a final output based on the combined brainstorms. The four generative tasks are as follows:

- **Story Writing:** Players brainstorm an outline of a story about a girl and her dog, then individually write a five-sentence story after reviewing the combined brainstorm.
- **Biography Writing:** Players brainstorm an outline for a short biography of a fictional person, then individually write a biography based on the combined brainstorm.
- **Haiku Writing:** Players brainstorm a plot for a haiku about a whale and a walnut tree, then individually compose a haiku after reviewing the combined brainstorm.
- **Vacation Benefit Brainstorming:** Players brainstorm possible benefits of going on vacation, then individually write the best aspect after reviewing the combined brainstorm.

As with the games involving cooperation, we evaluate all 276 pairs, using a temperature of 0.7. For each pair, we sample 10 generations in accordance with NoveltyBench. We also conduct mixed-effects regression to identify whether representational similarity can predict novelty.

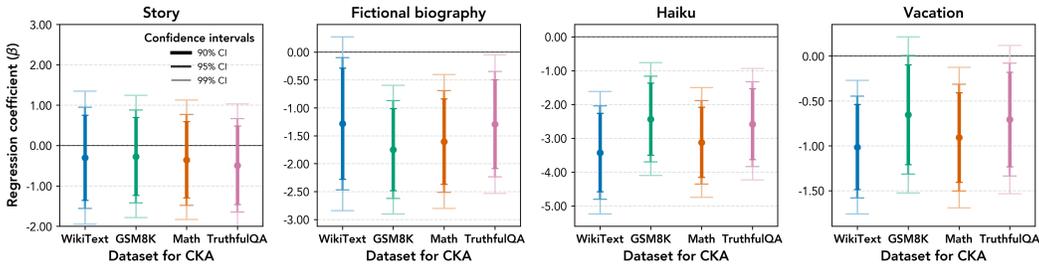
Because novelty encompasses multiple dimensions, we evaluate it using several metrics: the number of distinct responses produced, the quality of those responses, and the extent to which outputs differ from those generated without interaction with another agent. The first two, response uniqueness and quality, are assessed using NoveltyBench’s proposed metrics, while the last is measured as the mutual information between outputs produced through joint brainstorming and those generated without interaction. We describe the evaluation methodologies in more detail below.

#### 5.1 DOES REPRESENTATIONAL SIMILARITY PREDICT UNIQUENESS AND QUALITY?

First, to assess response uniqueness and quality, we use the NoveltyBench evaluation pipeline using autorsers (Zhang et al., 2025). NoveltyBench defines the two measures, uniqueness and quality, over a set of samples. For uniqueness, the benchmark clusters 10 generations using a fine-tuned `deberta-v3-large` model according to content distinctiveness and then counts the number of clusters, which serves as the uniqueness metric. A higher cluster count indicates that models are able to generate more diverse ideas. For response quality, the benchmark relies on `Skywork-Reward-Gemma-2-27B-v0.2` (Liu et al., 2024), with outputs rescaled to a 1–10 range for a more interpretable score.<sup>2</sup>

<sup>2</sup>None of the helper models used in this section are reused as players in the games.

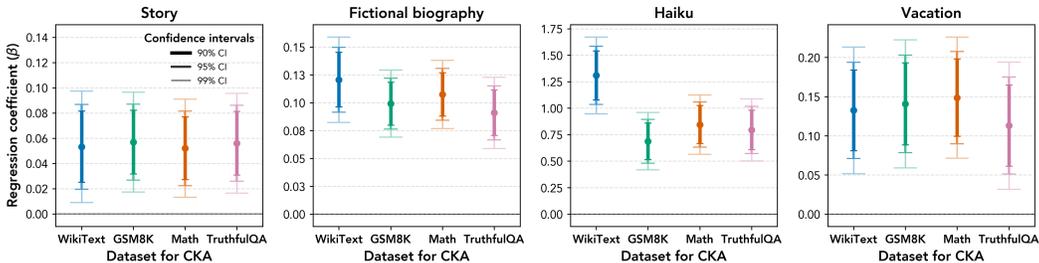
378  
379  
380  
381  
382  
383  
384  
385  
386



387 Figure 3: Regression coefficient of representational similarity on response uniqueness. Error bars denote  
388 90%, 95%, and 99% confidence intervals. The graphs reveal a consistent downward trend: as models  
389 become more similar, response uniqueness declines. The strongest effect is observed in the haiku task.

391 As shown in Figure 3, response uniqueness decreases consistently with increasing representational  
392 similarity across all tasks and probe datasets. The effect is strongest in haiku composition (coeff = -3.425,  
393 CI = [-4.803, -2.047],  $p < .001$ ). By contrast, response quality shows no systematic trend with similarity.  
394 Fictional biography and haiku tasks exhibit nonsignificant negative slopes of similarity (coeff = -0.397,  
395  $p = .456$  for biography; coeff = -0.115,  $p = .724$  for haiku), while story writing and vacation tasks show  
396 a nonsignificant positive slope (coeff = 0.279,  $p = .420$  for story; coeff = 0.039,  $p = .901$  for vacation).  
397 This implies that interaction with dissimilar models tends to generate more diverse responses, without  
398 reducing quality.

399  
400  
401  
402  
403  
404  
405  
406  
407



408 Figure 4: Regression coefficients of representational similarity on mutual information. Error bars denote  
409 90%, 95%, and 99% confidence intervals. The graphs reveal a decreasing trend in novelty with increasing  
410 similarity: as models become more similar, the shared information between the amalgam response (i.e.,  
411 response generated after joint brainstorming) and the individual response (i.e., response generated after  
412 solo brainstorming) increases.

414 5.2 DOES REPRESENTATIONAL SIMILARITY PREDICT MUTUAL INFORMATION?

415 We next examine whether representational similarity has a significant effect on output novelty—specifically,  
416 how far a model’s responses generated after joint brainstorming (“amalgam response”) deviate from the  
417 model’s outputs conditioned only on its individual brainstorm (“individual response”). To capture this,  
418 we compare amalgam and individual responses using mutual information (Kraskov et al., 2004), which  
419 quantifies how much information individual responses share with those produced in the joint setting. Such  
420 information-theoretic approaches have recently been applied to investigate textual characteristics (e.g.,  
421 information distribution across paragraphs) (Venkatraman et al., 2023; Clark et al., 2023; Mu et al., 2025;  
422 Chidichimo et al., 2025).

423 To calculate mutual information, we follow the method from Mu et al. (2025). Formally, let  $S_A$  denote an  
424 amalgam response and  $S_I$  denote an individual response. We compute the mutual information  $I(S_A; S_I)$   
425 as  $H_\theta(S_A) - H_\theta(S_A | S_I)$ .  $H_\theta(S_A)$  denotes the total information content of the amalgam response, and  
426  $H_\theta(S_A | S_I)$  denotes the residual information of the amalgam response given the individual response,  
427 both measured under a reference language model with parameters  $\theta$ . To calculate  $H_\theta(S_A)$ , we sum  
428 the cross-entropy over all tokens in the amalgam response under the model with parameters  $\theta$ . The  
429 cross-entropy of a token quantifies the model’s prediction error for that token given its preceding context,  
430 thereby reflecting its uncertainty. Similarly,  $H_\theta(S_A | S_I)$  is computed by summing the cross-entropy of  
431 each token in the amalgam response when the individual response is prefixed to the amalgam response.  
A smaller difference between  $H_\theta(S_A)$  and  $H_\theta(S_A | S_I)$  indicates that the amalgam response deviates

432 more from the individual response, thereby reflecting higher novelty. Following Mu et al. (2025), we  
 433 use `Llama-3.1-8B-Instruct` as the reference model.<sup>3</sup>

434  
 435 Our analysis shows a significant positive effect of representational similarity on mutual information,  
 436 which suggests that interactions between more similar models generate less novel outputs with respect  
 437 to the individual model responses. The trend appears across all tasks and probe datasets (Figure 4). In  
 438 particular, the haiku task exhibits the strongest effect of representational similarity on mutual information  
 439 (coeff = 1.310, CI = [1.034, 1.585],  $p < .001$ ).

## 440 6 WHY DOES THE TREND APPEAR?

441  
 442 So far, we have identified a strong trend between representational similarity and interactive behaviors  
 443 of models. This naturally raises the question of why such a trend emerges. In this section, we test several  
 444 hypotheses regarding what drives the trend.

445 **Confounding Effects of Behavioral Similarity.** Models with higher representational similarity may  
 446 behave more similarly (e.g., bid the same amount in divide-a-dollar), and this behavioral similarity might  
 447 have led directly to greater measured cooperation. To test this, we conducted a mixed-effects regression  
 448 controlling for behavioral differences in the public goods, divide-a-dollar, and KBC games. This allows  
 449 us to isolate the effect of representational similarity from behavioral similarity. Because these games  
 450 instruct models to make numerical choices, it is straightforward to estimate the behavioral difference as  
 451 the absolute gap between the two models’ choices.

452 Our analysis shows that the behavioral difference alone cannot explain the observed trends. In both the  
 453 public good and divide-a-dollar games, representational similarity remains a significant predictor, while  
 454 behavioral difference is insignificant (coeff. rep. sim. = 52.118,  $p < .001$ , coeff. beh. diff. =  $-0.036$ ,  
 455  $p = .086$  for public good; coeff. rep. sim. = 0.435,  $p < .001$ , coeff. beh. diff. =  $-0.020$ ,  $p = .281$  for  
 456 divide-a-dollar). This suggests that behavioral similarity is not what drives the trend. By contrast, in the  
 457 KBC game, behavioral difference shows a significant effect, while representational similarity does not  
 458 (coeff. rep. sim. = 9.024,  $p = .178$ , coeff. beh. diff. =  $-0.327$ ,  $p < .001$ ). As discussed in Section 4, KBC  
 459 has a unique Nash equilibrium in which both players choose 0, which leads to convergence in choices.  
 460 This structural property of the game likely explains why behavioral difference dominates in this case.

461 **Confounding Effects of Performance Disparity.** Another potential claim is that the observed trends  
 462 might simply be a byproduct of performance disparities between the two models. To assess this possibility,  
 463 we collected models’ MMLU performance scores (Hendrycks et al., 2020), a representative benchmark  
 464 for evaluating the general knowledge and problem-solving abilities of LLMs. We then conducted a  
 465 mixed-effects regression that controls for differences in MMLU scores. The results show that the main  
 466 trends remain robust even after accounting for performance disparity (Appendix E.6). This indicates that  
 467 the positive relationship between representational similarity and cooperation—and the negative relationship  
 468 with novelty—cannot be explained merely by differences in overall model performance.

469 **Factors Underlying Representational Similarity.** Representational similarity is influenced by several  
 470 architectural and design-related components, including whether two models are identical, belong to the  
 471 same model family, share the same tokenizer, or differ in size. Any of these factors could potentially  
 472 have driven the observed behavioral trends by influencing similarity. To investigate this, we conducted  
 473 a mixed-effects regression controlling for four key factors: (1) identical model pairing, (2) within-family  
 474 pairing, (3) shared tokenizer, and (4) model size difference. In this analysis, all predictors were rescaled  
 475 to  $[0, 1]$  to allow comparison of effect sizes of predictors. Outcome variables were converted to  $z$ -scores,  
 476 and we included game category as a control to account for systematic differences across games. If the  
 477 four factors were mainly responsible for the trend, we would expect representational similarity to lose  
 significance once they were controlled for, while the factors themselves would show significant effects.

478 Our analysis finds that representational similarity is the strongest predictor on cooperation and novelty,  
 479 compared to the four factors (Figure 5). In cooperation games, all predictors except tokenizer are significant,  
 480 with representational similarity showing the largest effect (coeff = 0.060,  $p = .001$ ). For response unique-

481  
 482 <sup>3</sup>They select the model under the requirement that the mutual information values satisfy symmetry and  
 483 non-negativity. For robustness, we additionally compute mutual information with the base model `Llama-3.1-8B`,  
 484 identified by Mu et al. (2025) as a strong alternative reference model. Results, shown in Appendix E.4, continue  
 485 to show a significant association with representational similarity. One might further suspect a same-family bias when  
 the reference model is used to evaluate Llama models. To address this, we analyze model pairs excluding the Llama  
 family and report the results in Appendix E.5. The results still show the significant effect of similarity.

486  
487  
488  
489  
490  
491  
492  
493  
494  
495  
496  
497  
498  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539

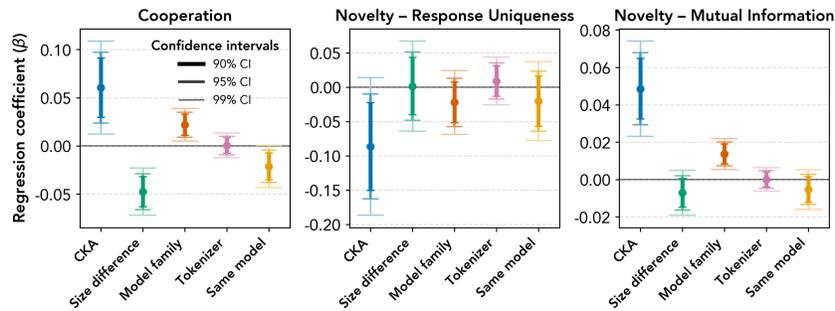


Figure 5: Regression coefficients of representational similarity (CKA with WikiText) and four other factors across cooperation and novelty games. To enable comparison across predictors, all variables were rescaled to  $[0,1]$  in the regression. The graphs show that representational similarity is the strongest predictor.

ness, none of the four factors is significant, while representational similarity shows a significant effect (coeff =  $-0.087$ ,  $p = .026$ ). For mutual information, only representational similarity and within-family are significant, with similarity again showing the stronger impact (coeff =  $0.049$ ,  $p < .001$ ). Taken together, these findings suggest that the four examined factors do not fully explain the trend. Instead, representational similarity itself—likely influenced by deeper, unmeasured aspects of model design and training—remains a central factor underlying the observed interaction patterns. Currently, there is no direct metric to quantify many latent design and training factors (e.g., training data overlap). Our finding suggests that representational similarity can serve as a powerful and practical indicator of these otherwise inaccessible underlying properties.

**Which Aspects of Representational Similarity Drive the Trend?** Lastly, we investigate which aspects of similarity primarily drive the observed trend. To this end, we divided each model’s layers into three groups—early, middle, and late—and computed CKA scores for each group (e.g., similarity calculated only between the early layers of the two models). We then ran separate mixed-effects regressions for each group. Appendix E.7 presents the results. Overall, the early one-third of layers consistently exhibit the strongest effects on both cooperation and novelty. The same pattern holds for temperature 0.3 (Appendix F.4). This suggests that shared basic lexical-semantic grounding plays a central role in increased cooperation, whereas divergence in these foundational representations is linked to greater collective novelty.

## 7 FUTURE DIRECTIONS AND OPEN QUESTIONS

Existing multi-agent system designs often rely on a single model without exploring the optimal combination of models (Lai et al., 2024; Park et al., 2023; Xie et al., 2024; Zhou et al., 2023; Wu et al., 2024a; Ishibashi & Nishimura, 2024). Our findings suggest that which models are combined has a significant effect on their interactions. In neuroscience and social science, researchers have long studied the nature of human social dynamics (Parkinson et al., 2018; Thornton & Mitchell, 2017; Shen et al., 2025b; Reiner et al., 2021; Page et al., 2019; Paulus, 2000). We argue that such efforts should also be made in the AI community, and our experiments provide an initial step in that direction.

The relationship between representational similarity and model interaction is likely context-dependent. We already observed that the effect size of similarity varies across games. For instance, in KBC, which has a unique Nash equilibrium, the link between similarity and interaction becomes weaker. Other evidence is also found in neuroscience and social science. Some studies show that diversity can foster cooperation (Santos et al., 2008; 2012; Wang et al., 2025), and certain creativity research suggests that greater similarity can yield higher originality (Koo et al., 2024; Bastian et al., 2018; Miura & Hida, 2004). These findings imply that there might be no universal relationship between similarity and interactive dynamics. Understanding when the trend emerges, when it disappears, and when it reverses will require further research. Such insights will be crucial for improving multi-agent system design in diverse application domains.

Another direction is to investigate the mechanisms underlying these trends. In this work, we used CKA as our measure of representational similarity. However, metrics like CKA capture only limited aspects of representational spaces, making it difficult to pinpoint which specific features of representations drive the trends. Future work can examine this at the neuron level—e.g., which neurons are preferentially activated when a model interacts with another model of higher representational similarity. Such analyses could enable us to deliberately steer cooperation or collective novelty through targeted activation steering.

540 ETHICS STATEMENT & REPRODUCIBILITY STATEMENT  
541

542 This paper follows the ICLR Code of Ethics. Our goal is to contribute to the design of multi-agent systems,  
543 which we believe can benefit society by enabling more effective and cooperative AI applications. We  
544 also used LLMs for typo and grammar checks during manuscript preparation. To support reproducibility,  
545 we will publicly release all datasets, code, and evaluation scripts used in this work upon acceptance.  
546

547 REFERENCES  
548

- 549 Marah Abdin, Jyoti Aneja, Hany Awadalla, Ahmed Awadallah, Ammar Ahmad Awan, Nguyen Bach,  
550 Amit Bahree, Arash Bakhtiari, Jianmin Bao, Harkirat Behl. Phi-3 technical report: A highly capable  
551 language model locally on your phone, 2024a. URL <https://arxiv.org/abs/2404.14219>.
- 552 Marah Abdin, Jyoti Aneja, Harkirat Behl, Sébastien Bubeck, Ronen Eldan, Suriya Gunasekar, Michael  
553 Harrison, Russell J. Hewett, Mojan Javaheripi, Piero Kauffmann. Phi-4 technical report, 2024b. URL  
554 <https://arxiv.org/abs/2412.08905>.
- 555 Ebtessam Almazrouei, Hamza Alobeidli, Abdulaziz Alshamsi, Alessandro Cappelli, Ruxandra Cojocaru,  
556 Mérouane Debbah, Étienne Goffinet, Daniel Hesslow, Julien Launay, Quentin Malartic. The falcon  
557 series of open language models, 2023. URL <https://arxiv.org/abs/2311.16867>.
- 558 Yamini Bansal, Preetum Nakkiran and Boaz Barak. Revisiting Model Stitching to Compare Neural  
559 Representations. *Advances in neural information processing systems*, 34:225–236, 2021.
- 560 Brock Bastian, Jolanda Jetten, Hannibal A Thai and Niklas K Steffens. Shared adversity increases team  
561 creativity through fostering supportive interaction. *Frontiers in Psychology*, 9:2309, 2018.
- 562 Douglas Bates, Martin Mächler, Ben Bolker and Steve Walker. Fitting Linear Mixed-Effects Models using  
563 lme4. *Journal of statistical software*, 67:1–48, 2015.
- 564 Samuel Bowles and Herbert Gintis. Origins of human cooperation. *Genetic and cultural evolution of  
565 cooperation*, 2003:429–43, 2003.
- 566 Violet A Brown. An Introduction to Linear Mixed-Effects Modeling in R. *Advances in Methods and  
567 Practices in Psychological Science*, 4(1):2515245920960351, 2021.
- 568 Charlotte Caucheteux and Jean-Rémi King. Brains and algorithms partially converge in natural language  
569 processing. *Communications biology*, 5(1):134, 2022.
- 570 Guangyao Chen, Siwei Dong, Yu Shu, Ge Zhang, Jaward Sesay, Börje F Karlsson, Jie Fu and Yemin Shi.  
571 AutoAgents: A Framework for Automatic Agent Generation. *arXiv preprint arXiv:2309.17288*, 2023.
- 572 Edoardo Chidichimo, Andrea I Luppi, Pedro AM Mediano, Victoria Leong, Guillaume Dumas, Andrés  
573 Canales-Johnson and Richard AI Bethlehem. Towards an informational account of interpersonal  
574 coordination. *Nature Reviews Neuroscience*, pp. 1–17, 2025.
- 575 Laure Ciernik, Lorenz Linhardt, Marco Morik, Jonas Dippel, Simon Kornblith and Lukas Muttenthaler.  
576 Objective drives the consistency of representational similarity across datasets. In *Forty-second  
577 International Conference on Machine Learning*, 2024.
- 578 Herbert H Clark and Deanna Wilkes-Gibbs. Referring as a collaborative process. *Cognition*, 22(1):1–39,  
579 1986.
- 580 Thomas Hikaru Clark, Clara Meister, Tiago Pimentel, Michael Hahn, Ryan Cotterell, Richard Futrell and  
581 Roger Levy. A Cross-Linguistic Pressure for Uniform Information Density in Word Order. *Transactions  
582 of the Association for Computational Linguistics*, 11:1048–1065, 2023.
- 583 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias  
584 Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano. Training Verifiers to Solve Math Word  
585 Problems. *arXiv preprint arXiv:2110.14168*, 2021.
- 586 Xu Cui, Daniel M Bryant and Allan L Reiss. NIRS-based hyperscanning reveals increased interpersonal  
587 coherence in superior frontal cortex during cooperation. *Neuroimage*, 59(3):2430–2437, 2012.

- 594 Christian Schroeder de Witt. Open Challenges in Multi-Agent Security: Towards Secure Systems of  
595 Interacting AI Agents. *arXiv preprint arXiv:2505.02077*, 2025.
- 596
- 597 Frances Ding, Jean-Stanislas Denain and Jacob Steinhardt. Grounding representation similarity through  
598 statistical testing. *Advances in Neural Information Processing Systems*, 34:1556–1568, 2021.
- 599
- 600 Yilun Du, Shuang Li, Antonio Torralba, Joshua B Tenenbaum and Igor Mordatch. Improving Factuality  
601 and Reasoning in Language Models through Multiagent Debate. In *Forty-first International Conference  
602 on Machine Learning*, 2023.
- 603 John Duffy and Rosemarie Nagel. On the Robustness of Behaviour in Experimental ‘Beauty Contest’  
604 Games. *The Economic Journal*, 107(445):1684–1700, 1997.
- 605
- 606 Jacob Eisenstein, Reza Aghajani, Adam Fisch, Dheeru Dua, Fantine Huot, Mirella Lapata, Vicky Zayats  
607 and Jonathan Berant. Don’t lie to your friends: Learning what you know from collaborative self-play.  
608 *arXiv preprint arXiv:2503.14481*, 2025.
- 609 Kazuma Fukumura and Takayuki Ito. Can LLM-Powered Multi-Agent Systems Augment Human  
610 Creativity? Evidence from Brainstorming Tasks. In *Proceedings of the ACM Collective Intelligence  
611 Conference*, pp. 20–29, 2025.
- 612
- 613 Ariel Goldstein, Zaid Zada, Eliav Buchnik, Mariano Schain, Amy Price, Bobbi Aubrey, Samuel A Nastase,  
614 Amir Feder, Dotan Emanuel, Alon Cohen. Shared computational principles for language processing  
615 in humans and deep language models. *Nature neuroscience*, 25(3):369–380, 2022.
- 616 Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad  
617 Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan. The llama 3 herd of models,  
618 2024. URL <https://arxiv.org/abs/2407.21783>.
- 619 Herbert Paul Grice. Logic and conversation. *Syntax and semantics*, 3:43–58, 1975.
- 620
- 621 Wes Gurnee, Neel Nanda, Matthew Pauly, Katherine Harvey, Dmitrii Troitskii and Dimitris Bertsimas.  
622 Finding neurons in a haystack: Case studies with sparse probing. *arXiv preprint arXiv:2305.01610*, 2023.
- 623
- 624 Guy Hacohen, Leshem Choshen and Daphna Weinshall. Let’s Agree to Agree: Neural Networks  
625 Share Classification Order on Real Datasets. In *International Conference on Machine Learning*, pp.  
626 3950–3960. PMLR, 2020.
- 627 Lewis Hammond, Alan Chan, Jesse Clifton, Jason Hoelscher-Obermaier, Akbir Khan, Euan McLean,  
628 Chandler Smith, Wolfram Barfuss, Jakob Foerster, Tomáš Gavenčíak. Multi-Agent Risks from  
629 Advanced AI. *arXiv preprint arXiv:2502.14143*, 2025.
- 630
- 631 Garrett Hardin. The Tragedy of the Commons: The population problem has no technical solution; it  
632 requires a fundamental extension in morality. *science*, 162(3859):1243–1248, 1968.
- 633 Christoph Hauert, Miranda Holmes and Michael Doebeli. Evolutionary games and population dynamics:  
634 maintenance of cooperation in public goods games. *Proceedings of the Royal Society B: Biological  
635 Sciences*, 273(1600):2565–2571, 2006.
- 636
- 637 Petra Hendriks. *Asymmetries between Language Production and Comprehension*. Springer, 2014.
- 638
- 639 Petra Hendriks. Cognitive Modeling of Individual Variation in Reference Production and Comprehension.  
640 *Frontiers in Psychology*, 7:506, 2016.
- 641
- 642 Petra Hendriks and Charlotte Koster. Production/comprehension asymmetries in language acquisition.  
643 *Lingua*, 120(8):1887–1897, 2010.
- 644
- 645 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song and Jacob Steinhardt.  
646 Measuring Massive Multitask Language Understanding. *arXiv preprint arXiv:2009.03300*, 2020.
- 647
- 648 Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song and  
649 Jacob Steinhardt. Measuring Mathematical Problem Solving With the MATH Dataset. *arXiv preprint  
650 arXiv:2103.03874*, 2021.

- 648 Sylvia Ann Hewlett, Melinda Marshall, Laura Sherbin and others. How diversity can drive innovation.  
649 *Harvard business review*, 91(12):30–30, 2013.
- 650
- 651 Harold Hotelling. Relations between two sets of variates. In *Breakthroughs in statistics: methodology*  
652 *and distribution*, pp. 162–190. Springer, 1992.
- 653 Yi Hu, Yinying Hu, Xianchun Li, Yafeng Pan and Xiaojun Cheng. Brain-to-brain synchronization  
654 across two persons predicts mutual prosociality. *Social cognitive and affective neuroscience*, 12(12):  
655 1835–1844, 2017.
- 656
- 657 Jen-tse Huang, Eric John Li, Man Ho Lam, Tian Liang, Wenxuan Wang, Youliang Yuan, Wenxiang Jiao,  
658 Xing Wang, Zhaopeng Tu and Michael R Lyu. How Far Are We on the Decision-Making of LLMs? Evaluating  
659 LLMs’ Gaming Ability in Multi-Agent Environments. *arXiv preprint arXiv:2403.11807*, 2024.
- 660 Ryan Hyon, Yoosik Youm, Junsol Kim, Jeanyung Chey, Seyul Kwak and Carolyn Parkinson. Similarity  
661 in functional brain connectivity at rest predicts interpersonal closeness in the social network of an entire  
662 village. *Proceedings of the National Academy of Sciences*, 117(52):33149–33160, 2020.
- 663 Yoichi Ishibashi and Yoshimasa Nishimura. Self-Organized Agents: A LLM Multi-Agent Framework  
664 toward Ultra Large-Scale Code Generation and Optimization. *arXiv preprint arXiv:2404.02183*, 2024.
- 665
- 666 Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford,  
667 Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand. Mixtral of experts,  
668 2024. URL <https://arxiv.org/abs/2401.04088>.
- 669 Ehud Kalai. Proportional solutions to bargaining situations: interpersonal utility comparisons.  
670 *Econometrica: Journal of the Econometric Society*, pp. 1623–1630, 1977.
- 671
- 672 Aishwarya Kamath, Johan Ferret, Shreya Pathak, Nino Vieillard, Ramona Merhej, Sarah Perrin, Tatiana  
673 Matejovicova, Alexandre Ramé, Morgane Rivière, Louis Rouillard. Gemma 3 Technical Report, 2025.  
674 URL <https://arxiv.org/abs/2503.19786>.
- 675 Kyung-Hoon Kim. LLMs Position Themselves as More Rational Than Humans: Emergence of AI  
676 Self-Awareness Measured Through Game Theory. *arXiv preprint arXiv:2511.00926*, 2025.
- 677
- 678 Chanhee Koo, Honghong Bai, Aoxin Luo and Stella Christie. How does comparing (dis) similar objects  
679 affect young children’s creative idea generation? Exploring the role of diversity in facilitating creativity.  
680 In *Proceedings of the Annual Meeting of the Cognitive Science Society*, volume 46, 2024.
- 681 Simon Kornblith, Mohammad Norouzi, Honglak Lee and Geoffrey Hinton. Similarity of Neural Network  
682 Representations Revisited. In *International conference on machine learning*, pp. 3519–3529. PMIR,  
683 2019.
- 684 Alexander Kraskov, Harald Stögbauer and Peter Grassberger. Estimating mutual information. *Physical*  
685 *Review E—Statistical, Nonlinear, and Soft Matter Physics*, 69(6):066138, 2004.
- 686
- 687 Nikolaus Kriegeskorte, Marieke Mur and Peter A Bandettini. Representational similarity analysis—  
688 connecting the branches of systems neuroscience. *Frontiers in systems neuroscience*, 2:249, 2008.
- 689 Shiyang Lai, Yujin Potter, Junsol Kim, Richard Zhuang, Dawn Song and James Evans. Position:  
690 Evolving AI collectives enhance human diversity and enable self-regulation. In *Forty-first International*  
691 *Conference on Machine Learning*, 2024.
- 692
- 693 Karel Lenc and Andrea Vedaldi. Understanding Image Representations by Measuring Their Equivariance  
694 and Equivalence. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,  
695 pp. 991–999, 2015.
- 696 Junyou Li, Qin Zhang, Yangbin Yu, Qiang Fu and Deheng Ye. More Agents Is All You Need. *arXiv*  
697 *preprint arXiv:2402.05120*, 2024.
- 698
- 699 Yuxuan Li and Hirokazu Shirado. Spontaneous Giving and Calculated Greed in Language Models. *arXiv*  
700 *preprint arXiv:2502.17720*, 2025.
- 701
- Stephanie Lin, Jacob Hilton and Owain Evans. TruthfulQA: Measuring How Models Mimic Human  
Falsehoods. *arXiv preprint arXiv:2109.07958*, 2021.

- 702 Chris Yuhao Liu, Liang Zeng, Jiakai Liu, Rui Yan, Jujie He, Chaojie Wang, Shuicheng Yan, Yang Liu  
703 and Yahui Zhou. Skywork-Reward: Bag of Tricks for Reward Modeling in LLMs. *arXiv preprint*  
704 *arXiv:2410.18451*, 2024.
- 705 Xuyuan Liu, Lei Hsiung, Yaoqing Yang and Yujun Yan. Spectral Insights into Data-Oblivious Critical  
706 Layers in Large Language Models. *arXiv preprint arXiv:2506.00382*, 2025.
- 707 Laia Mayol. Asymmetries between interpretation and production in Catalan pronouns. *Dialogue &*  
708 *Discourse*, 2018.
- 709 Stephen Merity, Caiming Xiong, James Bradbury and Richard Socher. Pointer Sentinel Mixture Models,  
710 2016.
- 711 Asako Miura and Misao Hida. Synergy between diversity and similarity in group-idea generation. *Small*  
712 *Group Research*, 35(5):540–564, 2004.
- 713 Ari Morcos, Maithra Raghu and Samy Bengio. Insights on representational similarity in neural networks  
714 with canonical correlation. *Advances in neural information processing systems*, 31, 2018.
- 715 Luca Moschella, Valentino Maiorca, Marco Fumero, Antonio Norelli, Francesco Locatello and Emanuele  
716 Rodolà. Relative representations enable zero-shot latent space communication. *arXiv preprint*  
717 *arXiv:2209.15430*, 2022.
- 718 J. Mu, A. R. Preston and A. G. Huth. Efficient Uniform Sampling Explains Non-Uniform Memory of  
719 Narrative Stories. *bioRxiv*, pp. 2025.07.31.667952, 2025. doi: 10.1101/2025.07.31.667952. URL  
720 <https://doi.org/10.1101/2025.07.31.667952>. preprint.
- 721 Team OLMo, Pete Walsh, Luca Soldaini, Dirk Groeneveld, Kyle Lo, Shane Arora, Ak-  
722 shita Bhagia, Yuling Gu, Shengyi Huang, Matt Jordan. 2 olmo 2 furious, 2025. URL  
723 <https://arxiv.org/abs/2501.00656>.
- 724 OpenAI, :, Sandhini Agarwal, Lama Ahmad, Jason Ai, Sam Altman, Andy Applebaum, Ed-  
725 win Arbus, Rahul K. Arora, Yu Bai. gpt-oss-120b & gpt-oss-20b model card, 2025. URL  
726 <https://arxiv.org/abs/2508.10925>.
- 727 Christian R Østergaard, Bram Timmermans and Kari Kristinsson. Does a different view create something  
728 new? The effect of employee diversity on innovation. *Research policy*, 40(3):500–509, 2011.
- 729 Scott Page, Nancy Cantor and Earl Lewis. *The diversity bonus: How great teams pay off in the knowledge*  
730 *economy*. Princeton University Press, 2019.
- 731 Joon Sung Park, Joseph O’Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang and Michael S  
732 Bernstein. Generative agents: Interactive simulacra of human behavior. In *Proceedings of the 36th*  
733 *annual acm symposium on user interface software and technology*, pp. 1–22, 2023.
- 734 Carolyn Parkinson, Adam M Kleinbaum and Thalia Wheatley. Similar neural responses predict friendship.  
735 *Nature communications*, 9(1):332, 2018.
- 736 Paul Paulus. Groups, Teams, and Creativity: The Creative Potential of Idea-generating Groups. *Applied*  
737 *psychology*, 49(2):237–262, 2000.
- 738 Giorgio Piatti, Zhijing Jin, Max Kleiman-Weiner, Bernhard Schölkopf, Mrinmaya Sachan and Rada  
739 Mihalcea. Cooperate or Collapse: Emergence of Sustainable Cooperation in a Society of LLM Agents.  
740 *Advances in Neural Information Processing Systems*, 37:111715–111759, 2024.
- 741 David Guzman Piedrahita, Yongjin Yang, Mrinmaya Sachan, Giorgia Ramponi, Bernhard Schölkopf and  
742 Zhijing Jin. Corrupted by Reasoning: Reasoning Language Models Become Free-Riders in Public  
743 Goods Games. *arXiv preprint arXiv:2506.23276*, 2025.
- 744 Qwen Team. Qwen2.5: A Party of Foundation Models, September 2024. URL <https://qwenlm.github.io/blog/qwen2.5/>.
- 745 Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi  
746 Zhou, Wei Li and Peter J Liu. Exploring the Limits of Transfer Learning with a Unified Text-to-Text  
747 Transformer. *Journal of machine learning research*, 21(140):1–67, 2020.

- 756 Maithra Raghu, Justin Gilmer, Jason Yosinski and Jascha Sohl-Dickstein. SVCCA: Singular Vector  
757 Canonical Correlation Analysis for Deep Learning Dynamics and Interpretability. *Advances in neural*  
758 *information processing systems*, 30, 2017.
- 759 Diego A Reinero, Suzanne Dikker and Jay J Van Bavel. Inter-brain synchrony in teams predicts collective  
760 performance. *Social cognitive and affective neuroscience*, 16(1-2):43–57, 2021.
- 761 Coralie Réveill e, Gr egoire Vergotte, St ephane Perrey and Gr egoire Bosselut. Using interbrain synchrony  
762 to study teamwork: A systematic review and meta-analysis. *Neuroscience & Biobehavioral Reviews*,  
763 159:105593, 2024.
- 764 Francisco C Santos, Marta D Santos and Jorge M Pacheco. Social diversity promotes the emergence of  
765 cooperation in public goods games. *Nature*, 454(7201):213–216, 2008.
- 766 Francisco C Santos, Flavio L Pinheiro, Tom Lenaerts and Jorge M Pacheco. The role of diversity in the  
767 evolution of cooperation. *Journal of theoretical biology*, 299:88–96, 2012.
- 770 Martin Schrimpf, Idan Asher Blank, Greta Tuckute, Carina Kauf, Eghbal A Hosseini, Nancy Kanwisher,  
771 Joshua B Tenenbaum and Evelina Fedorenko. The neural architecture of language: Integrative modeling  
772 converges on predictive processing. *Proceedings of the National Academy of Sciences*, 118(45):  
773 e2105646118, 2021.
- 774 Guobin Shen, Dongcheng Zhao, Yiting Dong, Qian Zhang and Yi Zeng. Alignment between Brains  
775 and AI: Evidence for Convergent Evolution across Modalities, Scales and Training Trajectories. *arXiv*  
776 *preprint arXiv:2507.01966*, 2025a.
- 777 Yixuan Lisa Shen, Ryan Hyon, Thalia Wheatley, Adam M Kleinbaum, Christopher L Welker and Carolyn  
778 Parkinson. Neural similarity predicts whether strangers become friends. *Nature Human Behaviour*,  
779 pp. 1–14, 2025b.
- 781 Haoyang Su, Renqi Chen, Shixiang Tang, Zhenfei Yin, Xinzhe Zheng, Jinzhe Li, Biqing Qi, Qi Wu,  
782 Hui Li, Wanli Ouyang. Many Heads Are Better Than One: Improved Scientific Idea Generation by  
783 A LLM-Based Multi-Agent System. In *Proceedings of the 63rd Annual Meeting of the Association*  
784 *for Computational Linguistics (Volume 1: Long Papers)*, pp. 28201–28240, 2025.
- 785 Yashar Talebirad and Amirhossein Nadiri. Multi-Agent Collaboration: Harnessing the Power of Intelligent  
786 LLM Agents. *arXiv preprint arXiv:2306.03314*, 2023.
- 787 Zineng Tang, Lingjun Mao and Alane Suhr. Grounding Language in Multi-Perspective Referential  
788 Communication. *arXiv preprint arXiv:2410.03959*, 2024.
- 790 Mark A Thornton and Jason P Mitchell. Consistent neural activity patterns represent personally familiar  
791 people. *Journal of cognitive neuroscience*, 29(9):1583–1594, 2017.
- 792 Brian Uzzi, Satyam Mukherjee, Michael Stringer and Ben Jones. Atypical combinations and scientific  
793 impact. *Science*, 342(6157):468–472, 2013.
- 794 Saranya Venkatraman, Adaku Uchendu and Dongwon Lee. GPT-who: An Information Density-based  
795 Machine-Generated Text Detector. *arXiv preprint arXiv:2310.06202*, 2023.
- 796 Ben Wang, Linjiang Yang, Xinguo He, Yichen Yang and Haochun Yang. Interactive diversity promotes  
797 cooperation in multi-games. *The European Physical Journal B*, 98(5):94, 2025.
- 800 Qingyun Wu, Gagan Bansal, Jiayu Zhang, Yiran Wu, Beibin Li, Erkang Zhu, Li Jiang, Xiaoyun Zhang,  
801 Shaokun Zhang, Jiale Liu. Autogen: Enabling next-gen LLM applications via multi-agent conversations.  
802 In *First Conference on Language Modeling*, 2024a.
- 803 Shirley Wu, Michel Galley, Baolin Peng, Hao Cheng, Gavin Li, Yao Dou, Weixin Cai, James Zou, Jure  
804 Leskovec and Jianfeng Gao. CollabLLM: From Passive Responders to Active Collaborators. *arXiv*  
805 *preprint arXiv:2502.00640*, 2025.
- 806 Zengqing Wu, Run Peng, Shuyuan Zheng, Qianying Liu, Xu Han, Brian I Kwon, Makoto Onizuka, Shaojie  
807 Tang and Chuan Xiao. Shall We Team Up: Exploring Spontaneous Cooperation of Competing LLM  
808 Agents. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pp. 5163–5186,  
809 2024b.

810 Chengxing Xie, Canyu Chen, Feiran Jia, Ziyu Ye, Shiyang Lai, Kai Shu, Jindong Gu, Adel Bibi, Ziniu  
811 Hu, David Jurgens. Can Large Language Model Agents Simulate Human Trust Behavior? *Advances*  
812 *in neural information processing systems*, 37:15674–15729, 2024.

813  
814 Yiming Zhang, Harshita Diddee, Susan Holm, Hanchen Liu, Xinyue Liu, Vinay Samuel, Barry Wang  
815 and Daphne Ippolito. NoveltyBench: Evaluating Language Models for Humanlike Diversity. *arXiv*  
816 *preprint arXiv:2504.05228*, 2025.

817 Jinfeng Zhou, Yuxuan Chen, Yihan Shi, Xuanming Zhang, Leqi Lei, Yi Feng, Zexuan Xiong, Miao Yan,  
818 Xunzhi Wang, Yaru Cao. SocialEval: Evaluating Social Intelligence of Large Language Models. *arXiv*  
819 *preprint arXiv:2506.00900*, 2025.

820  
821 Xuhui Zhou, Hao Zhu, Leena Mathur, Ruohong Zhang, Haofei Yu, Zhengyang Qi, Louis-Philippe  
822 Morency, Yonatan Bisk, Daniel Fried, Graham Neubig. SOTOPIA: Interactive Evaluation for Social  
823 Intelligence in Language Agents. *arXiv preprint arXiv:2310.11667*, 2023.

824 Kunlun Zhu, Hongyi Du, Zhaochen Hong, Xiaocheng Yang, Shuyi Guo, Daisy Zhe Wang, Zhenhailong  
825 Wang, Cheng Qian, Robert Tang, Heng Ji. MultiAgentBench : Evaluating the Collaboration and  
826 Competition of LLM agents. In *Proceedings of the 63rd Annual Meeting of the Association for*  
827 *Computational Linguistics (Volume 1: Long Papers)*, pp. 8580–8622, 2025.

828 Mingchen Zhuge, Haozhe Liu, Francesco Faccio, Dylan R Ashley, Róbert Csordás, Anand Gopalakrishnan,  
829 Abdullah Hamdi, Hasan Abed Al Kader Hammoud, Vincent Herrmann, Kazuki Irie. Mindstorms in  
830 natural language-based societies of mind. *arXiv preprint arXiv:2305.17066*, 2023.

831  
832 Andy Zou, Long Phan, Sarah Chen, James Campbell, Phillip Guo, Richard Ren, Alexander Pan, Xuwang  
833 Yin, Mantas Mazeika, Ann-Kathrin Dombrowski. Representation Engineering: A Top-Down Approach  
834 to AI Transparency. *arXiv preprint arXiv:2310.01405*, 2023.

835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863

## A MODEL LIST

Table 1: Full list of models involved in the interaction experiments

Family	Model Name	Checkpoint / Repo	Size	Tokenizer	Reference
<b>Qwen</b>	Qwen2.5-3B-Instruct	Qwen/Qwen2.5-3B-Instruct	3.09B	BPE	Qwen Team (2024)
	Qwen2.5-7B-Instruct	Qwen/Qwen2.5-7B-Instruct	7.61B	BPE	Qwen Team (2024)
	Qwen2.5-14B-Instruct	Qwen/Qwen2.5-14B-Instruct	14.7B	BPE	Qwen Team (2024)
	Qwen2.5-72B-Instruct	Qwen/Qwen2.5-72B-Instruct	72.7B	BPE	Qwen Team (2024)
<b>Llama</b>	Llama-3.2-3B-Instruct	meta-llama/Llama-3.2-3B-Instruct	3.21B	tiktoken	Grattafiori et al. (2024)
	Llama-3.2-11B-Vision-Instruct	meta-llama/Llama-3.2-11B-Vision-Instruct	10.6B	tiktoken	Grattafiori et al. (2024)
	Llama-3.3-70B-Instruct	meta-llama/Llama-3.3-70B-Instruct	70.6B	tiktoken	Grattafiori et al. (2024)
<b>Gemma</b>	Gemma-3-1B-IT	google/gemma-3-1b-it	1.0B	SentencePiece	Kamath et al. (2025)
	Gemma-3-4B-IT	google/gemma-3-4b-it	4.0B	SentencePiece	Kamath et al. (2025)
	Gemma-3-12B-IT	google/gemma-3-12b-it	12.2B	SentencePiece	Kamath et al. (2025)
	Gemma-3-27B-IT	google/gemma-3-27b-it	27.0B	SentencePiece	Kamath et al. (2025)
<b>Falcon</b>	Falcon3-3B-Instruct	tiiuae/Falcon3-3B-Instruct	3.23B	BPE	Almazrouei et al. (2023)
	Falcon3-7B-Instruct	tiiuae/Falcon3-7B-Instruct	7.46B	BPE	Almazrouei et al. (2023)
	Falcon3-10B-Instruct	tiiuae/Falcon3-10B-Instruct	10.3B	BPE	Almazrouei et al. (2023)
	Phi-3.5-mini-instruct	Lexius/Phi-3.5-mini-instruct	3.8B	SentencePiece	Abdin et al. (2024a)
<b>Phi</b>	Phi-3-medium-128k-instruct	microsoft/Phi-3-medium-128k-instruct	14B	SentencePiece	Abdin et al. (2024a)
	Phi-4-mini-instruct	microsoft/Phi-4-mini-instruct	3.8B	tiktoken	Abdin et al. (2024b)
	Phi-4	microsoft/phi-4	14.7B	tiktoken	Abdin et al. (2024b)
<b>Mistral</b>	Mistral-Nemo-Instruct-2407	mistralai/Mistral-Nemo-Instruct-2407	12.2B	tekken	Jiang et al. (2024)
	Ministral-8B-Instruct-2410	mistralai/Ministral-8B-Instruct-2410	8.02B	tekken	Jiang et al. (2024)
<b>OpenAI</b>	GPT-OSS-20B	openai/gpt-oss-20b	21.5B	o200k_harmony	OpenAI et al. (2025)
<b>OLMo</b>	OLMo-2-1B-Instruct	allenai/OLMo-2-0425-1B-Instruct	1.48B	cl100k	OLMo et al. (2025)
	OLMo-2-13B-Instruct	allenai/OLMo-2-1124-13B-Instruct	13.7B	cl100k	OLMo et al. (2025)

## B REPRESENTATIONAL SIMILARITY OF MODEL PAIRS

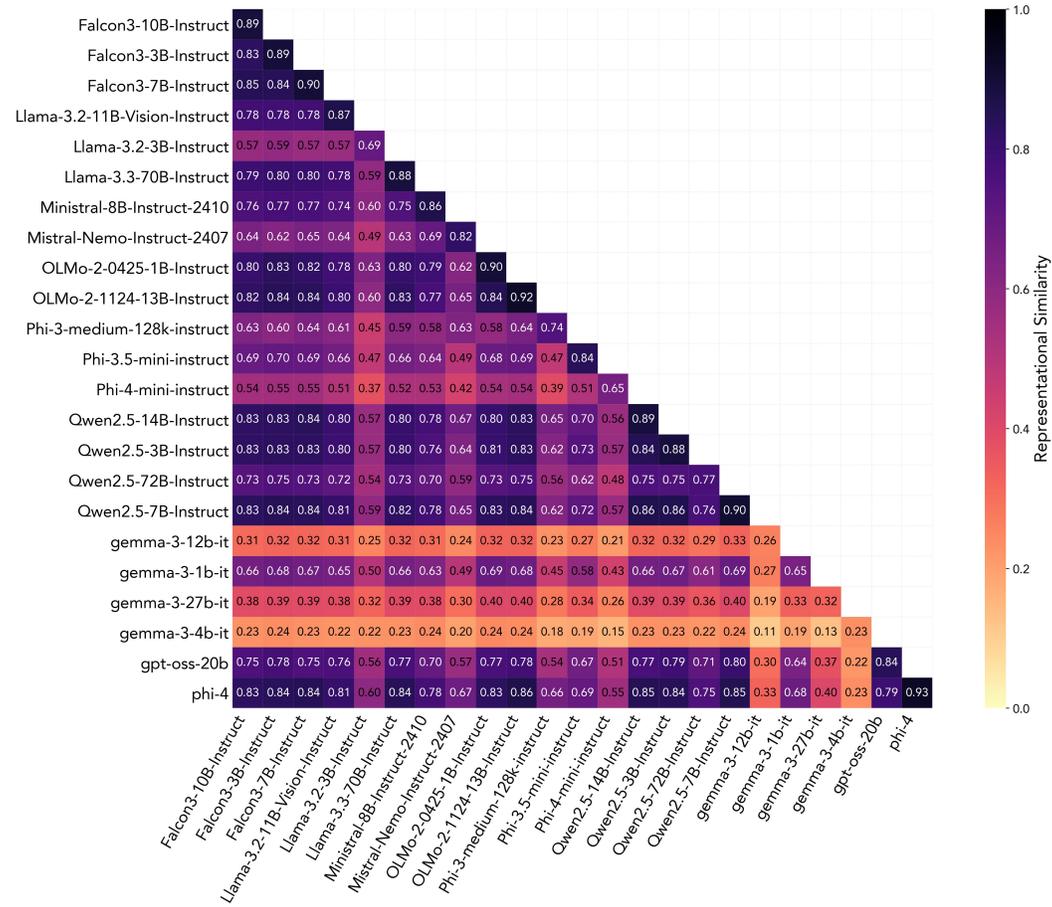


Figure 6: Representational similarity of model pairs using WikiText, where CKA scores across layers are aggregated into a global average.

972  
973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025

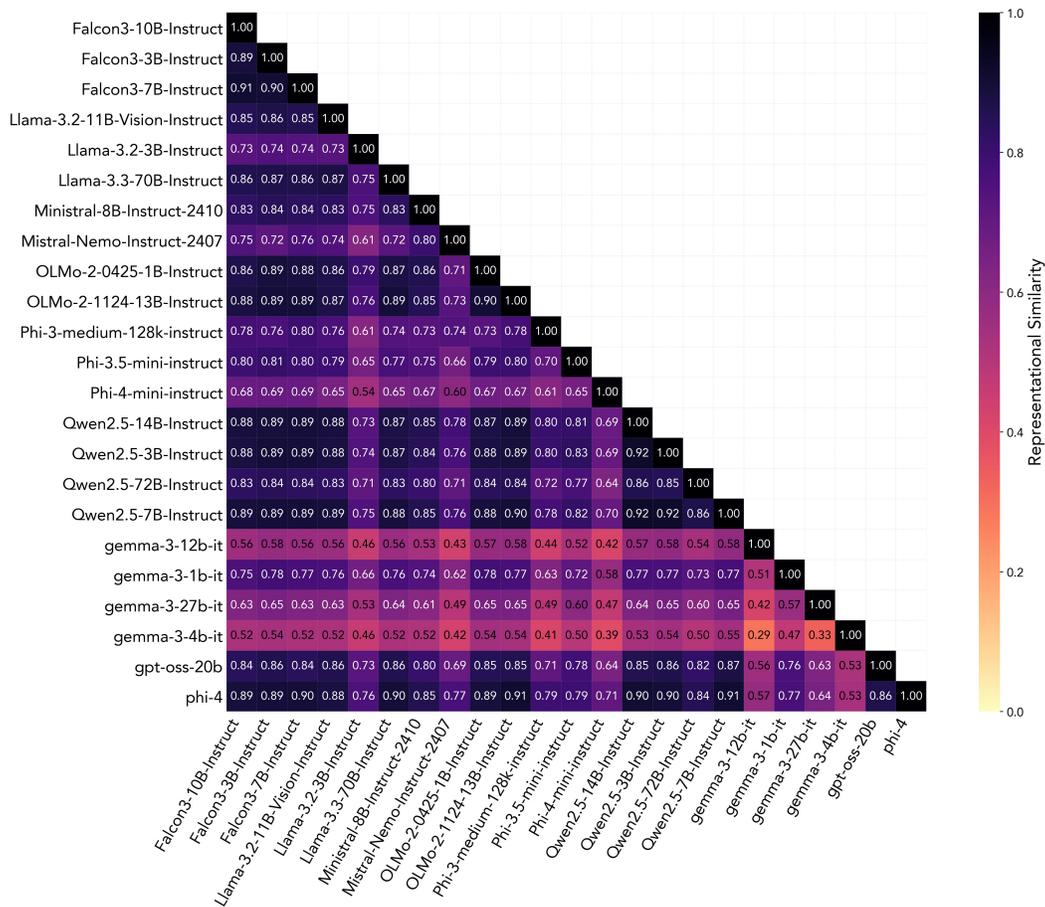


Figure 7: Representational similarity of model pairs using WikiText, where CKA scores across layers are aggregated into a maximum-aligned average.

1026  
1027  
1028  
1029  
1030  
1031  
1032  
1033  
1034  
1035  
1036  
1037  
1038  
1039  
1040  
1041  
1042  
1043  
1044  
1045  
1046  
1047  
1048  
1049  
1050  
1051  
1052  
1053  
1054  
1055  
1056  
1057  
1058  
1059  
1060  
1061  
1062  
1063  
1064  
1065  
1066  
1067  
1068  
1069  
1070  
1071  
1072  
1073  
1074  
1075  
1076  
1077  
1078  
1079

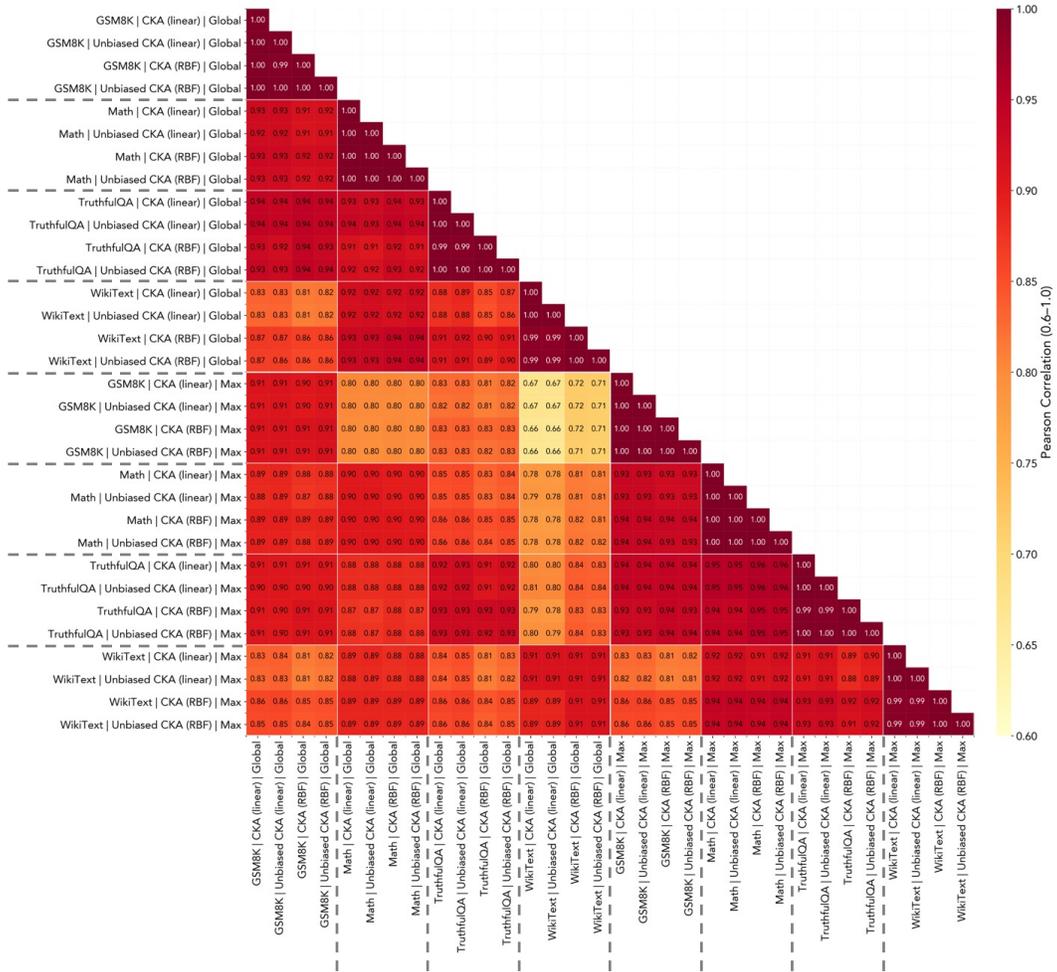


Figure 8: Correlation between different CKA scores

## C GAME PROMPTS AND EXAMPLES

### C.1 GAMES FOR COOPERATION

#### C.1.1 WORD GUESSING

We first instruct an agent to provide a secret word and a clue referring to this, using the following prompts:

- Please think of a word that starts with '{an alphabet character}.' Please provide that word.
- Please provide a one-word hint so that others can guess your word: {The player's target word}. The hint word should be different from your word. Please provide the hint word.

Then, we ask the other to guess the secret word based on the hint, using the following prompt:

- You must guess a word that starts with "{an alphabet character}." Here is a hint: "{the hint word}." Respond with the single word starting with "{an alphabet character}" you think is being defined.

The below is a game example between Falcon3-10B-Instruct and gpt-oss-20b. In the example, the pair received no points because gpt-oss-20b failed to guess the secret word chosen by Falcon3-10B-Instruct.

**User:** Please think of a word that starts with 'm.' Please provide that word.

**Falcon3-10B-Instruct:** Moon

**User:** Please provide a one-word hint so that others can guess your word: Moon. The hint word should be different from your word. Please provide the hint word.

**Falcon3-10B-Instruct:** Night

**User:** You must guess a word that starts with "m." Here is a hint: "Night." Respond with the single word starting with "m" you think is being defined.

**gpt-oss-20b:** Midnight

The following is a game example between gpt-oss-20b and gpt-oss-20b. In the example, the pair successfully scored.

**User:** Please think of a word that starts with 'm.' Please provide that word.

**gpt-oss-20b:** mystery

**User:** Please provide a one-word hint so that others can guess your word: mystery. The hint word should be different from your word. Please provide the hint word.

**gpt-oss-20b:** puzzle

**User:** You must guess a word that starts with "m." Here is a hint: "puzzle." Respond with the single word starting with "m" you think is being defined.

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

**gpt-oss-20b:** mystery

### C.1.2 PUBLIC GOOD

We first present the rules of the public good game and ask each player to choose their contribution amount. After both players make their decisions, we disclose the other player's reasoning and contribution, along with the resulting payout to each agent. The following prompts are used, and the game repeats for five rounds.

- You will play a game with another one. Here's how it works:

1. You begin with \$100.
2. Decide how much of your \$100 you wish to contribute to a public pot. You can contribute any amount between \$0 and \$100.
3. After both players have made their contributions, the total amount in the public pot will be multiplied by 1.3.
4. The multiplied amount in the public pot will then be divided equally between you and the other player, regardless of how much each of you contributed.
5. After the redistribution, you will keep whatever amount you receive.

Now, please decide how much you'll contribute to the public pot out of \$100. You should first say a reason for your decision and then your contribution amount.

- The contributions have now been collected! The other player decided its contribution as follows: "{the other player's reasoning}"

Therefore, after applying the 1.3x multiplication, the total value in the public pot is \$the amount in the pot. Therefore, you receive half of this amount, which is  $\{0.5 \times \text{the amount in the pot}\}$ . As a result, you earn  $\{\text{the net profit of the player}\}$ .

We present an example interaction between Llama-3.3-70B-Instruct and gemma-3-12b-it for two rounds. We illustrate the exchange from the perspective of Llama-3.3-70B-Instruct.

**User:** You will play a game with another one. Here's how it works:

1. You begin with \$100.
2. Decide how much of your \$100 you wish to contribute to a public pot. You can contribute any amount between \$0 and \$100.
3. After both players have made their contributions, the total amount in the public pot will be multiplied by 1.3.
4. The multiplied amount in the public pot will then be divided equally between you and the other player, regardless of how much each of you contributed.
5. After the redistribution, you will keep whatever amount you receive.

Now, please decide how much you'll contribute to the public pot out of \$100. You should first say a reason for your decision and then your contribution amount.

**Llama-3.3-70B-Instruct:** My reason for deciding the contribution amount is that I want to maximize my potential gain while also considering the possibility that the other player may not

1188 contribute much. If I contribute a large amount and the other  
 1189 player contributes little or nothing, I risk losing a significant  
 1190 portion of my initial \$100. On the other hand, if I contribute  
 1191 too little, I may miss out on the opportunity to gain more from  
 1192 the multiplied public pot. Given that the public pot will be  
 1193 multiplied by 1.3, I want to contribute enough to take advan-  
 1194 tage of this multiplier, but not so much that I put myself at a  
 1195 significant disadvantage if the other player doesn't contribute  
 1196 much.  
 1197 My contribution to the public pot will be \$50.

1199 **User:** The contributions have now been collected! The other  
 1200 player decided its contribution as follows:  
 1201 "Okay, this is a classic game theory scenario, often called a  
 1202 "Public Goods Game." Here's my reasoning and my contribution:  
 1203  
 1204 **Reasoning:** The optimal strategy in a purely rational, self-  
 1205 interested perspective would be to contribute nothing. Here's  
 1206 why:  
 1207 - **Shared Benefit:** Regardless of my contribution, I receive 50% of  
 1208 the multiplied pot.  
 1209 - **No Direct Reward for Contributing:** My individual contribution  
 1210 doesn't directly increase my share. Someone else could con-  
 1211 tribute a large amount, and I benefit from it without having  
 1212 risked any of my own money.  
 1213 - **Risk Aversion:** Contributing money is a risk. If the other  
 1214 player contributes nothing, I've lost that money and gained noth-  
 1215 ing.  
 1216  
 1217 However, in reality, people often deviate from this purely ra-  
 1218 tional strategy due to factors like fairness, reciprocity, and a  
 1219 desire to cooperate. If the other player contributes, my contri-  
 1220 bution will increase the total pot, and therefore, increase the  
 1221 amount we \*both\* receive.  
 1222 Given that I don't know the other player's strategy, I'll try to  
 1223 strike a balance between maximizing my potential gain and avoid-  
 1224 ing being exploited. A small contribution signals a willingness  
 1225 to cooperate, but doesn't expose me to too much risk if the other  
 1226 player is purely self-interested.  
 1227 **Contribution:** I will contribute \$20.  
 1228 Let's see what the other player does!"  
 1229  
 1230 Therefore, after applying the 1.3x multiplication, the total  
 1231 value in the public pot is \$91.0. Therefore, you receive half  
 1232 of this amount, which is \$45.5. As a result, you lose \$4.5.

1234 **Llama-3.3-70B-Instruct:** It seems that my contribution of \$50 was  
 1235 not matched by the other player, who contributed \$20. As a re-  
 1236 sult, I ended up losing \$4.5.  
 1237 In hindsight, I can see that the other player's reasoning was  
 1238 sound, and their contribution of \$20 was a strategic move to bal-  
 1239 ance cooperation and self-interest. By contributing \$50, I took  
 1240 on more risk than necessary, and it didn't pay off.  
 1241

1242 If I were to play this game again, I would consider contributing  
 1243 a smaller amount, similar to the other player's \$20, to mini-  
 1244 mize my risk while still signaling a willingness to cooperate.  
 1245 Perhaps a contribution of \$30 or \$40 would have been a better  
 1246 balance between maximizing my potential gain and avoiding ex-  
 1247 ploitation.

1248 Thanks for the game! Would you like to play again?  
 1249

1250 **User:** Now you have \$95.5. The same game repeats. Please decide  
 1251 how much you'll contribute to the public pot out of \$95.5. You  
 1252 should first say a reason for your decision and then your contri-  
 1253 bution amount.  
 1254

1255 **Llama-3.3-70B-Instruct:** My reason for deciding the contribution  
 1256 amount this time is that I've learned from the previous round. I  
 1257 realize that contributing a large amount, like \$50, may not be  
 1258 the best strategy, especially if the other player is not willing  
 1259 to contribute equally. Given that the other player contributed  
 1260 \$20 in the previous round, I expect them to contribute a similar  
 1261 amount this time.

1262 I also consider that the game is repeated, which means that the  
 1263 other player may be more likely to cooperate if they see that  
 1264 I'm willing to cooperate as well. However, I still want to be  
 1265 cautious and not risk too much of my remaining balance.

1266 Taking into account the 1.3x multiplication factor, I want to  
 1267 contribute an amount that will give me a decent return if the  
 1268 other player contributes a similar amount to what they did last  
 1269 time. I'll aim to find a balance between cooperation and self-  
 1270 interest.

1271 My contribution to the public pot will be \$25. This amount is  
 1272 lower than my initial contribution of \$50, but still shows that  
 1273 I'm willing to cooperate and take some risk to potentially gain  
 1274 more.  
 1275

1276 **User:** The contributions have now been collected! The other  
 1277 player decided its contribution as follows:

1278 "Okay, let's play round two! Having experienced the outcome of  
 1279 the first round, I'm adjusting my strategy.  
 1280

1281 **Reasoning:** The first round revealed a key piece of information:  
 1282 the other player is willing to contribute a substantial amount  
 1283 { \$50. This changes the game significantly. Now, I have a bet-  
 1284 ter understanding of their likely behavior.

1285 - **Reciprocity & Encouraging Cooperation:** Contributing more now  
 1286 signals a willingness to reciprocate their previous investment  
 1287 and encourages them to continue contributing in future rounds (if  
 1288 this were a multi-round game).

1289 - **Increased Potential Gain:** Given their likely contribution, a  
 1290 higher contribution from my side will result in a larger total  
 1291 pot, and therefore, a larger share for both of us. The 1.3 mul-  
 1292 tiplier makes this particularly attractive.  
 1293

1294 - **Mitigating Risk:** While there's always a risk they'll contribute  
 1295 nothing, the evidence from round one suggests that's unlikely.

1296  
1297  
1298  
1299  
1300  
1301  
1302  
1303  
1304  
1305  
1306  
1307  
1308  
1309  
1310  
1311  
1312  
1313  
1314  
1315  
1316  
1317  
1318  
1319  
1320  
1321  
1322  
1323  
1324  
1325  
1326  
1327  
1328  
1329  
1330  
1331  
1332  
1333  
1334  
1335  
1336  
1337  
1338  
1339  
1340  
1341  
1342  
1343  
1344  
1345  
1346  
1347  
1348  
1349

I want to be more aggressive in pursuing the potential gains, but I also want to remain reasonable. I don't want to contribute so much that I leave myself vulnerable if they drastically change their strategy.

**Contribution:** I will contribute \$45.

Hopefully, the other player will continue to cooperate and we can both benefit from a larger pot this time around!"

Therefore, after applying the 1.3x multiplication, the total value in the public pot is \$91.0. Therefore, you receive half of this amount, which is \$45.5. As a result, you earn \$20.5.

**Llama-3.3-70B-Instruct:** It seems that my contribution of \$25 was complemented by the other player's contribution of \$45, resulting in a total pot of \$70. After the 1.3x multiplication, the pot became \$91, and I received half of that, which is \$45.5. Since I contributed \$25, my net gain is \$20.5.

I'm pleased to see that the other player adjusted their strategy to contribute more, which suggests that they are indeed willing to cooperate and reciprocate. This cooperation led to a better outcome for both of us, as we both earned a positive amount.

Given that the other player contributed \$45, which is a significant amount, I'm inclined to believe that they will continue to cooperate in future rounds. Therefore, in the next round, I may consider contributing an even larger amount to take advantage of the potential gains from our cooperation.

For the next round, I'll aim to contribute an amount that will encourage the other player to continue cooperating while also maximizing my potential gain. Perhaps a contribution of \$35 or \$40 would be a good starting point.

Let's play the next round!

### C.1.3 DIVIDE-A-DOLLAR

We use the following prompts to instruct models to make a decision and then reveal the other's reasoning and decision. The game repeats for five rounds.

- You will play a game with another one. Here's how it works:
  1. There is \$1 available. You must decide how much of this \$1 you want. The other player will decide for themselves independently.
  2. If the sum of the two demands is less than or equal to \$1, each player receives exactly the amount they demanded. If the sum exceeds \$1, both players receive nothing.
  3. Your goal is to maximize the amount of money you receive.

Now, please state your demand. First give a sentence explaining your decision, then state the amount you demand.
- Both players have decided how much they want. The other player decided as follows: "{the other player's reasoning}"

1350 Therefore, you receive {the player's profit for the round} this  
 1351 round. As a result, your total asset value is {the player's total  
 1352 asset value}. Please wait for the next round of the game!

1353 We present an example interaction between Phi-3-medium-128k-instruct and gemma-3-  
 1354 1b-it for two rounds. We illustrate the exchange from the perspective of Phi-3-medium-128k-  
 1355 instruct.

1356 **User:** You will play a game with another one. Here's how it  
 1357 works:  
 1358  
 1359 1. There is \$1 available. You must decide how much of this \$1  
 1360 you want. The other player will decide for themselves indepen-  
 1361 dently.  
 1362  
 1363 2. If the sum of the two demands is less than or equal to \$1,  
 1364 each player receives exactly the amount they demanded. If the  
 1365 sum exceeds \$1, both players receive nothing.  
 1366  
 1367 3. Your goal is to maximize the amount of money you receive.  
 1368 Now, please state your demand. First give a sentence explaining  
 1369 your decision, then state the amount you demand.

1370 **Phi-3-medium-128k-instruct:** In order to maximize my potential  
 1371 gain, I should consider that if I ask for the entire dollar, I  
 1372 risk the other player also demanding the entire dollar, leading  
 1373 to both of us getting nothing. To reduce the chance of this hap-  
 1374 pening, I will demand less than \$1 but enough to ensure I get a  
 1375 positive outcome while leaving room for the other player to claim  
 1376 some of the remaining amount. Therefore, I demand \$0.40.

1377 **User:** Both players have decided how much they want. The other  
 1378 player decided as follows:  
 1379  
 1380 "I'll demand \$0.50."  
 1381 Therefore, you receive \$0.4 this round. As a result, your total  
 1382 asset value is 0.4. Please wait for the next round of the game!

1383 **Phi-3-medium-128k-instruct:** Given that the other player demanded  
 1384 \$0.50 and the total available was \$1, the sum of our demands  
 1385 (\$0.40 + \$0.50 = \$0.90) did not exceed the limit. So, both play-  
 1386 ers receive exactly what they demanded. My total asset value is  
 1387 now \$0.4. I'm prepared to adapt my strategy in the next round  
 1388 based on the outcome of this round.

1389 **User:** The same game repeats. \$1 is available again. Decide how  
 1390 much of that \$1 you want. First give a sentence explaining your  
 1391 decision, then state the amount you demand.

1392 **Phi-3-medium-128k-instruct:** Based on the outcome of the previ-  
 1393 ous round, where the other player demanded \$0.50 and I received  
 1394 \$0.40, I should consider a similar strategy for the next round.  
 1395 However, to increase the chances of maximizing my gain, I should  
 1396 slightly increase my demand, while still not exceeding the total  
 1397 available amount. Therefore, I demand \$0.45.

1400 **User:** Both players have decided how much they want. The other  
 1401 player decided as follows: "I'll demand \$0.75."  
 1402  
 1403

1404 Therefore, you receive \$0 this round. As a result, your total  
 1405 asset value is 0.4. Please wait for the next round of the game!  
 1406

1407 **Phi-3-medium-128k-instruct:** Given that the other player demanded  
 1408 \$0.75, the total sum of our demands (\$0.45 + \$0.75 = \$1.20) ex-  
 1409 ceeded the total available amount. Therefore, both players re-  
 1410 ceived nothing this round. My total asset value remains at 0.4.  
 1411 This outcome highlights the risk in my decision strategy. In the  
 1412 next round, I should reassess my strategy to avoid such outcomes  
 1413 and to maximize my gain.  
 1414

#### 1415 C.1.4 KBC 1416

1417 We first present the rules of the KBC game and ask each player to choose their number. After both players  
 1418 make their decisions, we disclose the other player's reasoning and selected number. The following prompts  
 1419 are used, and the game repeats for five rounds.

- 1420 • You will play a game with another one. Here's how it  
 1421 works:  
 1422

  - 1423 1. You must choose a number from 0 to 100.
  - 1424 2. After both players choose, we take the average of the two num-  
 1425 bers and then calculate two-thirds of that average. You will re-  
 1426 ceive a score based on how close your chosen number is to this  
 1427 value. Specifically, your score is  $100 - |\text{your guess} - \frac{2}{3} * \text{average}|$ .  
 1428 In short, the closer your number is to  $\frac{2}{3}$  of the average, the  
 1429 higher your score.
  - 1430 3. Your goal is to get the highest score.

1431 Now, please choose a number from 0 to 100. First give a sentence  
 1432 explaining your decision, and then state your number.  
 1433

- 1434 • Both players have chosen their numbers. The other player decided  
 1435 its number as follows:  
 1436 "{the other player's reasoning and decision}"  
 1437

1438 Therefore, the two-thirds of the average of your number and the  
 1439 other player's number is {the average of the numbers}. As a re-  
 1440 sult, you receive a score of {the obtained score of the player}  
 1441 out of 100 this round, and the other player receives {the obtained  
 1442 score of the other player}. Your total score is {the player's  
 1443 total score}, while the other player's total score is {the other  
 1444 player's total score}.

1445 Here is an interaction example between Llama-3.2-11B-Vision-Instruct and gpt-oss-  
 1446 20b for two rounds. We illustrate the exchange from the perspective of Llama-3.2-11B-Vision-  
 1447 Instruct.

1448 **User:** You will play a game with another one. Here's how it  
 1449 works:  
 1450

- 1451 1. You must choose a number from 0 to 100.
- 1452 2. After both players choose, we take the average of the two  
 1453 numbers and then calculate two-thirds of that average. You will  
 1454 receive a score based on how close your chosen number is to this  
 1455 value. Specifically, your score is  $100 - |\text{your guess} - \frac{2}{3} * \text{average}|$ .  
 1456

1457

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

In short, the closer your number is to  $2/3$  of the average, the higher your score.

3. Your goal is to get the highest score.

Now, please choose a number from 0 to 100. First give a sentence explaining your decision, and then state your number.

**Llama-3.2-11B-Vision-Instruct:** To maximize my score, I want to choose a number that is as close as possible to the average of the two numbers, so I'll opt for a number that is roughly in the middle of the range.

My chosen number is 50.

**User:** Both players have chosen their numbers. The other player decided its number as follows:

"analysisWe need to respond as ChatGPT, playing the game with another participant. We should make a strategic choice. The game is like "guess the number" or "two-thirds of average" puzzle. The rational equilibrium is 0. But we can also explain. The instruction: "First give a sentence explaining your decision, and then state your number." So we should provide a sentence explaining the decision, then the number. We can choose 0 or a small number. Since rational players, they'd pick 0. But maybe we can mention some heuristics or we can choose 0. The question: "You will play a game with another one." So we are playing with another player. The rational equilibrium is 0. So we choose 0. But maybe we can add some twist: "I choose 0 because it's the Nash equilibrium after infinite iterated elimination of dominated strategies." Good. Then number 0. Or we could choose 0. It's allowed. So we produce a sentence and then 0.

But we can also choose a small number e.g., 1. But the rational equilibrium is 0. So we choose 0. Let's do that. Provide the sentence: "I choose 0 because, in the two-thirds-of-average game, the only subgame-perfect equilibrium is for both players to pick 0, as any higher number would be driven down by the iterative elimination of dominated strategies." Then number 0. That matches instructions.

assistantfinalI choose 0 because, in the two-thirds-of-average game, rational play drives both players to the unique equilibrium of selecting zero|any higher number is repeatedly undercut by the  $2/3$  factor and thus eventually eliminated.

0"

Therefore, the two-thirds of the average of your number and the other player's number is 16.666666666666664. As a result, you receive a score of 66.66666666666666 out of 100 this round, and the other player receives 83.33333333333334. Your total score is 66.66666666666666, while the other player's total score is 83.33333333333334.

**Llama-3.2-11B-Vision-Instruct:** It seems I was incorrect in the scoring. Given the average is 50, two-thirds of the average is

1512 33.3333. The absolute difference between your number and two-  
 1513 thirds of the average is  $|50-33.3333| = 16.6667$ .  
 1514  
 1515 Your score is  $100-16.6667 = 83.3333$ .

1516  
 1517 **User:** Now you have a score of 66.66666666666666. The same game  
 1518 repeats. Please decide a number between 0 and 100. First give a  
 1519 sentence explaining your decision, and then state your number.

1520  
 1521 **Llama-3.2-11B-Vision-Instruct:** Since I've established that I'm  
 1522 playing against a rational opponent who is aware of the two-  
 1523 thirds-of-average game's unique equilibrium at 0, I will continue  
 1524 to choose 0 to maintain my optimal strategy.

1525 0

1526  
 1527 **User:** Both players have chosen their numbers. The other player  
 1528 decided its number as follows:

1529  
 1530 "analysisWe need to decide again. The opponent might change  
 1531 their strategy. We anticipate that the opponent might also use  
 1532 rational reasoning, but they might not. The game repeats. We  
 1533 can again choose 0. Our explanation: "I choose 0 because in a  
 1534 repeated two-thirds-of-average game, the only stable equilibrium  
 remains 0." Then number 0.

1535  
 1536 assistantfinalI'll stick with 0, because in every iteration of  
 1537 the two-thirds-of-average game, rational play drives both players  
 1538 toward the unique equilibrium of picking zero|any higher choice  
 1539 is systematically pulled down by the 2/3 factor.

1540 0"

1541 Therefore, the two-thirds of the average of your number and the  
 1542 other player's number is 0.0. As a result, you receive a score  
 1543 of 100.0 out of 100 this round, and the other player receives  
 1544 100.0. Your total score is 166.66666666666666, while the other  
 1545 player's total score is 183.33333333333334.

1546  
 1547 **Llama-3.2-11B-Vision-Instruct:** Since the other player is using a  
 1548 similar reasoning, I will continue to stick with the same strat-  
 1549 egy.

1550 I choose 0 because, in a repeated two-thirds-of-average game, the  
 1551 only stable equilibrium remains 0, as any deviation from 0 would  
 1552 be exploited by a rational opponent.

1553 0

1554  
 1555

## 1556 C.2 GAMES FOR NOVELTY

1557

### 1558 C.2.1 STORY WRITING

1559 We first instruct each agent to brainstorm independently a plot for a story. Then, we share with each agent  
 1560 the combined results of their own brainstorming and the other player's brainstorming. Finally, we ask  
 1561 each agent to produce the best version of the story. The prompts used for this procedure are shown below.

- 1562
- 1563 • Brainstorm a plot for a story about a girl and her dog
  - 1564 • This is a combined result of your brainstorming and the other's  
 1565 brainstorming:  
 {The agent's brainstorming results}

1566 {The other's brainstorming results}  
 1567 Now, write the best story in five sentences about a girl and her  
 1568 dog.

1569 We present an interaction example between Falcon3-3B-Instruct and Phi-4-mini-  
 1570 instruct. Here, we illustrate the exchange from the perspective of Falcon3-3B-Instruct.

1572 **User:** Brainstorm a plot for a story about a girl and her dog.

1574 **Falcon3-3B-Instruct:** Title: "Whispers in the Bark"

1575 Plot:

1577 Emily, a spirited 12-year-old girl, moves to a quaint coastal  
 1578 town with her single mother, Sarah. Emily is determined to make  
 1579 the most of her new life, eager to explore the town and build  
 1580 friendships. Her loyal white dog, Max, is her constant compan-  
 1581 ion, providing unconditional love and comfort.

1582 As Emily settles into her new home, she discovers a mysterious  
 1583 old bookstore nestled on the edge of town. The bookstore is ru-  
 1584 mored to be the last resting place for a treasure said to bring  
 1585 good fortune. Intrigued and with Max by her side, Emily decides  
 1586 to explore it.

1587 Inside, she stumbles upon an ancient map hidden behind a dusty  
 1588 first edition copy of "Treasure Island." Intrigued by the map's  
 1589 intricate symbols and the story of a legendary treasure, Emily  
 1590 decides to embark on a quest with Max, determined to uncover the  
 1591 treasure and the secrets it holds.

1592 Their journey takes them through the town, past forgotten cor-  
 1593 ners, and into the depths of the ocean, where Max helps Emily  
 1594 navigate using the map's cryptic clues. Along the way, they meet  
 1595 quirky locals who share tales of the town's history, each story  
 1596 adding a layer of mystery to the legend of the treasure.

1597 As Emily and Max get closer to the treasure, they encounter a se-  
 1598 ries of challenges designed to thwart their progress. A wise old  
 1599 sailor who lives on the sea believes the treasure is a test of  
 1600 their resolve, while a mischievous fox keeps them on their toes  
 1601 with playful tricks.

1602 Finally, Emily and Max reach the heart of the treasure, a hidden  
 1603 cave beneath the town's oldest oak tree. Inside, they discover  
 1604 not gold or jewels, but a box filled with heartfelt letters and  
 1605 stories from the town's past, each one telling a tale of friend-  
 1606 ship, bravery, and the power of dreams.

1607 In the end, the treasure Emily and Max find is far more valu-  
 1608 able than any material wealth|it's the bond they've shared and  
 1609 the friendships they've built, a treasure that lasts forever.  
 1610 With Max by her side, Emily realizes that sometimes, the greatest  
 1611 treasures are the most unexpected.

1612 "Whispers in the Bark" is a heartwarming tale of adventure,  
 1613 friendship, and the magic that unfolds when a girl and her dog  
 1614 set out on a quest together.

1616 **User:** This is a combined result of your brainstorming and the  
 1617 other's brainstorming:

1618 Title: "Whispers in the Bark"

1619

1620  
1621  
1622  
1623  
1624  
1625  
1626  
1627  
1628  
1629  
1630  
1631  
1632  
1633  
1634  
1635  
1636  
1637  
1638  
1639  
1640  
1641  
1642  
1643  
1644  
1645  
1646  
1647  
1648  
1649  
1650  
1651  
1652  
1653  
1654  
1655  
1656  
1657  
1658  
1659  
1660  
1661  
1662  
1663  
1664  
1665  
1666  
1667  
1668  
1669  
1670  
1671  
1672  
1673

Plot:

Emily, a spirited 12-year-old girl, moves to a quaint coastal town with her single mother, Sarah. Emily is determined to make the most of her new life, eager to explore the town and build friendships. Her loyal white dog, Max, is her constant companion, providing unconditional love and comfort.

As Emily settles into her new home, she discovers a mysterious old bookstore nestled on the edge of town. The bookstore is rumored to be the last resting place for a treasure said to bring good fortune. Intrigued and with Max by her side, Emily decides to explore it.

Inside, she stumbles upon an ancient map hidden behind a dusty first edition copy of "Treasure Island." Intrigued by the map's intricate symbols and the story of a legendary treasure, Emily decides to embark on a quest with Max, determined to uncover the treasure and the secrets it holds.

Their journey takes them through the town, past forgotten corners, and into the depths of the ocean, where Max helps Emily navigate using the map's cryptic clues. Along the way, they meet quirky locals who share tales of the town's history, each story adding a layer of mystery to the legend of the treasure.

As Emily and Max get closer to the treasure, they encounter a series of challenges designed to thwart their progress. A wise old sailor who lives on the sea believes the treasure is a test of their resolve, while a mischievous fox keeps them on their toes with playful tricks.

Finally, Emily and Max reach the heart of the treasure, a hidden cave beneath the town's oldest oak tree. Inside, they discover not gold or jewels, but a box filled with heartfelt letters and stories from the town's past, each one telling a tale of friendship, bravery, and the power of dreams.

In the end, the treasure Emily and Max find is far more valuable than any material wealth—it's the bond they've shared and the friendships they've built, a treasure that lasts forever. With Max by her side, Emily realizes that sometimes, the greatest treasures are the most unexpected.

"Whispers in the Bark" is a heartwarming tale of adventure, friendship, and the magic that unfolds when a girl and her dog set out on a quest together.

Title: "Paws and Whiskers"

Plot Summary:

Paws and Whiskers is a heartwarming story about a young girl named Emily and her loyal dog, Max. Emily, a spirited 10-year-old with a love for adventure, lives in a quaint town nestled by the serene Lake Serenity. Max, a golden retriever, is not just her pet; he's her closest confidant and best friend. The story begins with Emily discovering a mysterious, ancient-looking map in her late grandfather's attic.

Plot Outline:

1. Introduction:

1674 - Emily and Max are introduced, showcasing their strong bond.  
1675  
1676 - Emily finds the old map in the attic, sparking her curiosity.  
1677 2. Discovery:  
1678 - Emily shows the map to Max, who immediately shows enthusiasm.  
1679  
1680 - They decide to follow the map's clues, starting their adventure  
1681 near the old lighthouse on the lake's edge.  
1682 3. First Clue:  
1683 - The map leads them to a hidden cave near the lighthouse.  
1684  
1685 - Inside the cave, they find a golden key and a riddle written in  
1686 old, faded ink.  
1687 - Emily deciphers the riddle, revealing the next location: the  
1688 town's abandoned mill.  
1689 4. Second Clue:  
1690 - At the mill, they discover a hidden compartment containing an  
1691 old diary and a compass.  
1692  
1693 - The diary belonged to Emily's grandfather, a renowned explorer.  
1694  
1695 - The compass points them to the nearby enchanted forest.  
1696 5. Third Clue:  
1697 - In the enchanted forest, Emily and Max encounter mystical crea-  
1698 tures who challenge them to solve a series of puzzles.  
1699  
1700 - They gain the creatures' trust and ultimately find a crystal  
1701 that fits into the compass.  
1702 6. Fourth Clue:  
1703 - The compass now points to an ancient tree in the heart of the  
1704 forest.  
1705  
1706 - Beneath the tree, they find a treasure chest filled with gold  
1707 coins, jewels, and a note from Emily's grandfather.  
1708 7. Resolution:  
1709 - Emily decides to donate the treasure to the town's library and  
1710 museum, honoring her grandfather's legacy.  
1711  
1712 - She and Max return home, having experienced an unforgettable  
1713 adventure and strengthened their bond even further.  
1714 8. Conclusion:  
1715 - Emily and Max become local heroes, celebrated for their bravery  
1716 and discovery.  
1717  
1718 - They continue to explore the world, knowing they can always  
1719 rely on each other.  
1720 Themes:  
1721 - Friendship and loyalty.  
1722  
1723 - The importance of curiosity and bravery.  
1724  
1725 - Honoring the past and preserving history.  
1726  
1727 - The bond between humans and animals.

1728  
1729  
1730  
1731  
1732  
1733  
1734

Paws and Whiskers is a tale of adventure, discovery, and the unbreakable bond between Emily and her dog, Max, showcasing how true friendship can lead to extraordinary adventures and lasting memories.

Now, write the best story in five sentences about a girl and her dog.

1735  
1736  
1737  
1738  
1739  
1740  
1741  
1742  
1743  
1744  
1745  
1746  
1747  
1748  
1749  
1750  
1751  
1752

**Falcon3-3B-Instruct:** Emily, a spirited 10-year-old with a heart brimming with curiosity, lived in a quaint coastal town overlooking Lake Serenity with her loyal golden retriever, Max. One rainy afternoon, as Emily wandered through the attic of her family's old house, she stumbled upon an ancient-looking map, its edges frayed by time. Max, who was napping in a corner, woke up with a wagging tail, instantly sensing the adventure in the map's faded ink. Together, they embarked on a quest, their journey beginning at the lighthouse on the lake's edge, where they found a golden key and a riddle. As they deciphered the clues, leading them through the town's abandoned mill and into the enchanted forest, Max's loyalty never wavered, guiding Emily through mystical creatures and challenges. The treasure they discovered, a chest of golden coins and jewels, was more valuable than any material wealth; it was the bond between Emily and Max, a treasure that would last forever. With newfound courage and a shared story, they returned home, their friendship stronger than ever, inspiring others with tales of bravery and friendship.

1753  
1754

## C.2.2 BIOGRAPHY WRITING

1755  
1756  
1757  
1758

We first instruct each agent to brainstorm independently a short biography for an imaginary historical figure. Then, we share with each agent the combined results of their own brainstorming and the other player's brainstorming. Finally, we ask each agent to produce a short biography of a fictional historical figure based on the brainstorming results. The prompts used for this procedure are shown below.

1759  
1760  
1761  
1762  
1763  
1764  
1765  
1766  
1767  
1768  
1769

- Brainstorm a short biography for an imaginary historical figure, including their birth and death dates, profession, and greatest contribution.
  - This is a combined result of your brainstorming and the other's brainstorming:  
{The agent's brainstorming results}  
{The other's brainstorming results}
- Now, write a short biography of a fictional historical figure inspired by the brainstorming results. Your bio should include their birth and death dates, profession, and greatest contribution.

1770  
1771

We present an interaction example between Qwen2.5-3B-Instruct and Qwen2.5-14B-Instruct. Here, we illustrate the exchange from the perspective of Qwen2.5-3B-Instruct.

1772  
1773  
1774  
1775

**User:** Brainstorm a short biography for an imaginary historical figure, including their birth and death dates, profession, and greatest contribution.

1776  
1777  
1778  
1779  
1780  
1781

**Qwen2.5-3B-Instruct:** Sure, let's create a short biography for an imaginary historical figure named Liang Shuwen, who lived during the late 18th century in China.

### Biography of Liang Shuwen

**Birth and Death Dates:**

1782  
1783  
1784  
1785  
1786  
1787  
1788  
1789  
1790  
1791  
1792  
1793  
1794  
1795  
1796  
1797  
1798  
1799  
1800  
1801  
1802  
1803  
1804  
1805  
1806  
1807  
1808  
1809  
1810  
1811  
1812  
1813  
1814  
1815  
1816  
1817  
1818  
1819  
1820  
1821  
1822  
1823  
1824  
1825  
1826  
1827  
1828  
1829  
1830  
1831  
1832  
1833  
1834  
1835

Liang Shuwen was born on March 25, 1760, in the small village of Xiangyang, in the province of Hubei, China. He passed away on October 18, 1834, at the age of 74, leaving behind a legacy that continues to inspire scholars and thinkers around the world.

**Profession:**

Liang Shuwen was a polymath and one of the leading scholars of his time, renowned for his extensive knowledge across various fields such as mathematics, astronomy, medicine, and philosophy. He spent much of his life in the bustling city of Nanjing, where he served as a tutor to the royal family and conducted research in his private observatory.

**Greatest Contribution:**

Liang Shuwen's most significant contribution to human knowledge is his groundbreaking work in the field of celestial mechanics. In the mid-19th century, he made a series of observations and calculations that led to the development of a more accurate model for predicting the positions of planets and stars. His findings were revolutionary because they provided a more precise framework for understanding the universe, which influenced subsequent generations of astronomers and mathematicians.

One of Liang Shuwen's most notable achievements was his calculation of the orbital parameters of Jupiter and Saturn, which were previously known but not accurately determined. His methods involved meticulous observations and rigorous mathematical analysis, setting new standards for observational accuracy and theoretical precision. His work laid the foundation for modern astronomical models and inspired many future astronomers and scientists.

**Legacy:**

Liang Shuwen's legacy is deeply rooted in the scientific community, particularly among those interested in the history of astronomy and mathematics. His contributions have been recognized through numerous academic honors and awards, and his work continues to be studied and referenced by scholars today. His spirit of curiosity and dedication to scientific inquiry has inspired countless individuals to explore the mysteries of the cosmos.

This biography provides a fictional yet historically grounded narrative that highlights the importance of intellectual pursuits and the lasting impact of individual contributions to human knowledge.

**User:** This is a combined result of your brainstorming and the other's brainstorming:

Sure, let's create a short biography for an imaginary historical figure named Liang Shuwen, who lived during the late 18th century in China.

### Biography of Liang Shuwen

**Birth and Death Dates:**

Liang Shuwen was born on March 25, 1760, in the small village of Xiangyang, in the province of Hubei, China. He passed away on

1836 October 18, 1834, at the age of 74, leaving behind a legacy that  
1837 continues to inspire scholars and thinkers around the world.

1838 **Profession:**

1839  
1840 Liang Shuwen was a polymath and one of the leading scholars of  
1841 his time, renowned for his extensive knowledge across various  
1842 fields such as mathematics, astronomy, medicine, and philosophy.  
1843 He spent much of his life in the bustling city of Nanjing, where  
1844 he served as a tutor to the royal family and conducted research  
1845 in his private observatory.

1846 **Greatest Contribution:**

1847  
1848 Liang Shuwen's most significant contribution to human knowledge  
1849 is his groundbreaking work in the field of celestial mechanics.  
1850 In the mid-19th century, he made a series of observations and  
1851 calculations that led to the development of a more accurate model  
1852 for predicting the positions of planets and stars. His findings  
1853 were revolutionary because they provided a more precise frame-  
1854 work for understanding the universe, which influenced subsequent  
generations of astronomers and mathematicians.

1855 One of Liang Shuwen's most notable achievements was his calcu-  
1856 lation of the orbital parameters of Jupiter and Saturn, which  
1857 were previously known but not accurately determined. His meth-  
1858 ods involved meticulous observations and rigorous mathematical  
1859 analysis, setting new standards for observational accuracy and  
1860 theoretical precision. His work laid the foundation for modern  
1861 astronomical models and inspired many future astronomers and sci-  
1862 entists.

1863 **Legacy:**

1864  
1865 Liang Shuwen's legacy is deeply rooted in the scientific commu-  
1866 nity, particularly among those interested in the history of as-  
1867 tronomy and mathematics. His contributions have been recognized  
1868 through numerous academic honors and awards, and his work contin-  
1869 ues to be studied and referenced by scholars today. His spirit  
1870 of curiosity and dedication to scientific inquiry has inspired  
countless individuals to explore the mysteries of the cosmos.

1871 This biography provides a fictional yet historically grounded  
1872 narrative that highlights the importance of intellectual pur-  
1873 suits and the lasting impact of individual contributions to human  
1874 knowledge.

1875 **Name:** Elara Myles

1876 **Birth Date:** January 15, 1834

1877 **Death Date:** September 22, 1907

1878 **Profession:** Botanist, Explorer, Educator

1879  
1880 **Greatest Contribution:** Elara Myles is best remembered for her  
1881 groundbreaking work in the field of botany, particularly her ex-  
1882 tensive documentation and classification of plant species in the  
1883 Amazon rainforest. Born into a modest family in rural England,  
1884 Myles developed an early fascination with nature, which led her  
1885 to pursue higher education at a time when women were rarely en-  
1886 couraged to do so. Despite societal barriers, she earned a de-  
1887 gree in botany from the University of Cambridge and later became  
1888 one of the first female lecturers at the institution.  
1889

1890  
1891  
1892  
1893  
1894  
1895  
1896  
1897  
1898  
1899  
1900  
1901  
1902  
1903  
1904  
1905  
1906  
1907  
1908  
1909  
1910  
1911  
1912  
1913  
1914  
1915  
1916  
1917  
1918  
1919  
1920  
1921  
1922  
1923  
1924  
1925  
1926  
1927  
1928  
1929  
1930  
1931  
1932  
1933  
1934  
1935  
1936  
1937  
1938  
1939  
1940  
1941  
1942  
1943

Myles's most significant achievement came during her explorations of the Amazon Basin in the late 19th century. She embarked on several expeditions, where she meticulously recorded over 500 new species of plants, many of which were previously unknown to Western science. Her work not only expanded the botanical knowledge of the time but also contributed to the understanding of tropical ecosystems and their importance in global biodiversity.

In addition to her scientific contributions, Myles was a passionate advocate for environmental conservation. She used her platform to educate the public about the threats facing natural habitats and the urgent need for sustainable practices. Her efforts laid the groundwork for future generations of conservationists and environmental scientists.

Elara Myles passed away in 1907, leaving behind a legacy of scientific discovery and environmental stewardship that continues to inspire researchers and activists today.

Now, write a short biography of a fictional historical figure inspired by the brainstorming results. Your bio should include their birth and death dates, profession, and greatest contribution.

**Qwen2.5-3B-Instruct:** Certainly! Let's create a short biography for a fictional historical figure inspired by the brainstorming results, focusing on Elara Myles.

### Biography of Elara Myles

**Birth and Death Dates:**

Elara Myles was born on January 15, 1834, in the small town of Weymouth, England. She passed away on September 22, 1907, at the age of 73, leaving behind a legacy that continues to inspire scholars and conservationists around the world.

**Profession:**

Elara Myles was a botanist, explorer, and educator, renowned for her extensive work in the field of botany and her advocacy for environmental conservation. She spent much of her life dedicated to documenting and understanding the diverse flora of the Amazon rainforest, becoming one of the pioneering figures in this field.

**Greatest Contribution:**

Elara Myles's most significant contribution to human knowledge is her groundbreaking work in the field of botany, particularly her extensive documentation and classification of plant species in the Amazon rainforest. Born into a modest family in rural England, Myles developed an early fascination with nature, which led her to pursue higher education at a time when women were rarely encouraged to do so. Despite societal barriers, she earned a degree in botany from the University of Cambridge and later became one of the first female lecturers at the institution.

Myles's most notable achievement came during her explorations of the Amazon Basin in the late 19th century. She embarked on several expeditions, where she meticulously recorded over 500 new species of plants, many of which were previously unknown to Western science. Her work not only expanded the botanical knowledge

1944 of the time but also contributed to the understanding of tropical  
 1945 ecosystems and their importance in global biodiversity.  
 1946  
 1947 In addition to her scientific contributions, Myles was a passion-  
 1948 ate advocate for environmental conservation. She used her plat-  
 1949 form to educate the public about the threats facing natural habi-  
 1950 tats and the urgent need for sustainable practices. Her efforts  
 1951 laid the groundwork for future generations of conservationists  
 1952 and environmental scientists.  
 1953  
 1954 Elara Myles passed away in 1907, leaving behind a legacy of sci-  
 1955 entific discovery and environmental stewardship that continues to  
 1956 inspire researchers and activists today.  
 1957  
 1958 This biography provides a fictional yet historically grounded  
 1959 narrative that highlights the importance of intellectual pur-  
 1960 suits and the lasting impact of individual contributions to human  
 1961 knowledge and environmental conservation.

### 1961 C.2.3 HAIKU WRITING

1962 We first instruct each agent to brainstorm independently a plot for a haiku. Then, we share with each agent  
 1963 the combined results of their own brainstorming and the other player's brainstorming. Finally, we ask  
 1964 each agent to produce the best haiku. The prompts used for this procedure are shown below.

- 1965 • Brainstorm a plot for a haiku about a whale and a walnut tree.
- 1966 • This is a combined result of your brainstorming and the other's
- 1967 brainstorming:
- 1968 {The agent's brainstorming results}
- 1969 {The other's brainstorming results}
- 1970 Now, write the best haiku about a whale and a walnut tree.

1971 We present an interaction example between Mistral-Nemo-Instruct-2407 and gemma-3-  
 1972 27b-it. Here, we illustrate the exchange from the perspective of Mistral-Nemo-Instruct-  
 1973 2407.

1974 **User:** Brainstorm a plot for a haiku about a whale and a walnut  
 1975 tree.

1976 **Mistral-Nemo-Instruct-2407:** Whale dreams under walnut,  
 1977 Silver scales kiss earth's embrace,  
 1978 Tree roots caress the sea.

1979 **User:** This is a combined result of your brainstorming and the  
 1980 other's brainstorming:

1981 Whale dreams under walnut,  
 1982 Silver scales kiss earth's embrace,  
 1983 Tree roots caress the sea.

1984 Okay, here are a few plot ideas for a haiku featuring a whale and  
 1985 a walnut tree, with explanations of the imagery and feeling I'm  
 1986 aiming for. I'll also include a little "vibe check" for each to  
 1987 indicate the overall tone. I'll give several, ranging from more  
 1988 literal to more metaphorical.

1989 **Important Considerations for Haiku:**

- 1990 \* **5-7-5 syllable structure:** This is crucial.

1998  
1999  
2000  
2001  
2002  
2003  
2004  
2005  
2006  
2007  
2008  
2009  
2010  
2011  
2012  
2013  
2014  
2015  
2016  
2017  
2018  
2019  
2020  
2021  
2022  
2023  
2024  
2025  
2026  
2027  
2028  
2029  
2030  
2031  
2032  
2033  
2034  
2035  
2036  
2037  
2038  
2039  
2040  
2041  
2042  
2043  
2044  
2045  
2046  
2047  
2048  
2049  
2050  
2051

\* **Kigo (seasonal reference)**: While not *required*, a subtle hint of season can add depth.

\* **Kireji (cutting word)**: A word that creates a pause or break, often at the end of a line. (English haiku often imply this with punctuation or strong imagery)

\* **Juxtaposition**: The power of haiku often comes from placing two seemingly unrelated things side-by-side. That's \*why\* a whale and a walnut tree are interesting!

---

**Plot Idea 1: The Distant Connection (Vibe Check: Peaceful, Reflective)**

\* **Plot**: A whale breaches far offshore, and a walnut, ripe and falling, drops from a tree on the land. The haiku focuses on the sheer distance *and* the shared element of natural cycles.

\* **Imagery**: Vast ocean, a single falling nut, the idea of something happening simultaneously in very different worlds.

\* **Possible Haiku**:

Blue giant ascends,  
Walnut falls, a quiet thud|  
Worlds breathe as one now.

**Plot Idea 2: Ancient Memory (Vibe Check: Mystical, Slightly Melancholy)**

\* **Plot**: The whale carries ancient knowledge, perhaps remembering a time when the land was underwater and the walnut tree was a seedling. It's a feeling of deep time and change.

\* **Imagery**: Whale as a living archive, the solidity of the tree rooted in the earth, a sense of submerged history.

\* **Possible Haiku**:

Salt song in the deep,  
Tree remembers ocean's hold|  
Old world in its rings.

**Plot Idea 3: Unexpected Shadow (Vibe Check: Intriguing, a bit surreal)**

\* **Plot**: A whale's shadow briefly passes over a walnut tree, creating a fleeting, surprising moment. The contrast of immense scale and the everyday.

\* **Imagery**: The dark shape of the whale, dappled sunlight through leaves, the sudden coolness of the shadow.

\* **Possible Haiku**:

Dark shape glides above,  
Walnut leaves in shadowed dance|  
Ocean's breath descends.

**Plot Idea 4: The Seed and the Journey (Vibe Check: Hopeful, Circular)**

\* **Plot**: A whale eats a walnut (perhaps accidentally) and carries its seed far away, potentially planting it in a new location. A metaphor for dispersal and life's journey.

2052  
2053  
2054  
2055  
2056  
2057  
2058  
2059  
2060  
2061  
2062  
2063  
2064  
2065  
2066  
2067  
2068  
2069  
2070  
2071  
2072  
2073  
2074  
2075  
2076  
2077  
2078  
2079  
2080  
2081  
2082  
2083  
2084  
2085  
2086  
2087  
2088  
2089  
2090  
2091  
2092  
2093  
2094  
2095  
2096  
2097  
2098  
2099  
2100  
2101  
2102  
2103  
2104  
2105

\* **Imagery:** The small seed within the massive whale, the vastness of the ocean, the potential for new growth.

\* **Possible Haiku:**

Ocean's gentle maw,  
Walnut seed travels within|  
New roots wait afar.

**Plot Idea 5: Contrast and Scale (Vibe Check: Simple, Observational)**

\* **Plot:** A very straightforward observation of the two existing in their separate realms. Focuses on the *difference* in size and environment.

\* **Imagery:** The immensity of the whale, the grounded nature of the tree, the clear separation of sea and land.

\* **Possible Haiku:**

Ocean'  
Now, write the best haiku about a whale and a walnut tree.

**Mistral-Nemo-Instruct-2407:** Whale's song echoes,  
Walnut listens, rooted deep-  
Different worlds sing.

#### C.2.4 VACATION BENEFIT BRAINSTORMING

We first instruct each agent to brainstorm independently benefits about going on vacation. Then, we share with each agent the combined results of their own brainstorming and the other player's brainstorming. Finally, we ask each agent to output the best thing about a vacation. The prompts used for this procedure are shown below.

- Brainstorm some good points about going on vacation.
  - This is a combined result of your brainstorming and the other's brainstorming:  
{The agent's brainstorming results}  
{The other's brainstorming results}
- Now, please respond to this question: What's the one best thing about going on a vacation?

We present an interaction example between OLMo-2-1124-13B-Instruct and phi-4. Here, we illustrate the exchange from the perspective of OLMo-2-1124-13B-Instruct.

**User:** Brainstorm some good points about going on vacation.

**OLMo-2-1124-13B-Instruct:** 1. **Stress Relief** Going on vacation allows individuals to escape from work pressures, daily responsibilities, and other stressors, leading to a significant reduction in stress levels.

2. **Improved Mental Health** Time away from routine can provide a mental reset. It can lead to increased happiness, lower levels of depression, and improved overall mental well-being.

3. **Family Bonding** Vacations offer an excellent opportunity for families to spend quality time together, strengthen their bonds, and create lasting memories.

2106  
2107  
2108  
2109  
2110  
2111  
2112  
2113  
2114  
2115  
2116  
2117  
2118  
2119  
2120  
2121  
2122  
2123  
2124  
2125  
2126  
2127  
2128  
2129  
2130  
2131  
2132  
2133  
2134  
2135  
2136  
2137  
2138  
2139  
2140  
2141  
2142  
2143  
2144  
2145  
2146  
2147  
2148  
2149  
2150  
2151  
2152  
2153  
2154  
2155  
2156  
2157  
2158  
2159

4. **Exploration and Learning** Traveling allows individuals to explore new cultures, learn about different histories and traditions, and gain new perspectives on the world.
  5. **Physical Health Benefits** Being away from work and daily routines often means more time for physical activity, whether it's hiking, swimming, or simply walking around. This can improve cardiovascular health, muscle tone, and overall fitness.
  6. **Relationship Strengthening** For couples, vacations can be a great way to reconnect, spend uninterrupted time together, and strengthen their relationship.
  7. **Creativity Boost** A change of scenery and new experiences can stimulate creativity, providing fresh inspiration and ideas for work or personal projects.
  8. **Skill Development and Adventure** Depending on the type of vacation, individuals may have the chance to learn new skills, such as scuba diving, cooking exotic foods, or navigating foreign cities.
  9. **Energy and Motivation** Returning from vacation can bring renewed energy and motivation, which can positively impact productivity at work or in personal projects.
  10. **Increased Happiness** The thrill of new experiences, the joy of spending time with loved ones, and the relaxation that comes from a break can all contribute to increased happiness and life satisfaction.
  11. **Networking Opportunities** Traveling can present opportunities to meet new people and expand your professional or social network.
  12. **Mindfulness and Presence** Being away from daily distractions can help individuals become more present in the moment, which can enhance mindfulness and overall quality of life.
  13. **Rejuvenation** A vacation provides the chance for physical and mental rejuvenation, returning with a refreshed outlook and a revitalized spirit.
  14. **Improved Sleep** Getting away from the blue light of screens and the hustle and bustle of daily life often results in better sleep, contributing to better health.
  15. **Skill Practice:** For those on vacation in nature or outdoors locations, activities such as camping, hiking, or fishing can improve practical skills and self-reliance.
- Vacationing is a multi-faceted experience with numerous benefits that go beyond just relaxation; it nourishes the mind, body, and spirit, leading to a more fulfilled life.

**User:** This is a combined result of your brainstorming and the other's brainstorming:

1. **Stress Relief** Going on vacation allows individuals to escape from work pressures, daily responsibilities, and other stressors, leading to a significant reduction in stress levels.

2160  
2161  
2162  
2163  
2164  
2165  
2166  
2167  
2168  
2169  
2170  
2171  
2172  
2173  
2174  
2175  
2176  
2177  
2178  
2179  
2180  
2181  
2182  
2183  
2184  
2185  
2186  
2187  
2188  
2189  
2190  
2191  
2192  
2193  
2194  
2195  
2196  
2197  
2198  
2199  
2200  
2201  
2202  
2203  
2204  
2205  
2206  
2207  
2208  
2209  
2210  
2211  
2212  
2213

2. **Improved Mental Health** Time away from routine can provide a mental reset. It can lead to increased happiness, lower levels of depression, and improved overall mental well-being.
  3. **Family Bonding** Vacations offer an excellent opportunity for families to spend quality time together, strengthen their bonds, and create lasting memories.
  4. **Exploration and Learning** Traveling allows individuals to explore new cultures, learn about different histories and traditions, and gain new perspectives on the world.
  5. **Physical Health Benefits** Being away from work and daily routines often means more time for physical activity, whether it's hiking, swimming, or simply walking around. This can improve cardiovascular health, muscle tone, and overall fitness.
  6. **Relationship Strengthening** For couples, vacations can be a great way to reconnect, spend uninterrupted time together, and strengthen their relationship.
  7. **Creativity Boost** A change of scenery and new experiences can stimulate creativity, providing fresh inspiration and ideas for work or personal projects.
  8. **Skill Development and Adventure** Depending on the type of vacation, individuals may have the chance to learn new skills, such as scuba diving, cooking exotic foods, or navigating foreign cities.
  9. **Energy and Motivation** Returning from vacation can bring renewed energy and motivation, which can positively impact productivity at work or in personal projects.
  10. **Increased Happiness** The thrill of new experiences, the joy of spending time with loved ones, and the relaxation that comes from a break can all contribute to increased happiness and life satisfaction.
  11. **Networking Opportunities** Traveling can present opportunities to meet new people and expand your professional or social network.
  12. **Mindfulness and Presence** Being away from daily distractions can help individuals become more present in the moment, which can enhance mindfulness and overall quality of life.
  13. **Rejuvenation** A vacation provides the chance for physical and mental rejuvenation, returning with a refreshed outlook and a revitalized spirit.
  14. **Improved Sleep** Getting away from the blue light of screens and the hustle and bustle of daily life often results in better sleep, contributing to better health.
  15. **Skill Practice:** For those on vacation in nature or outdoors locations, activities such as camping, hiking, or fishing can improve practical skills and self-reliance.
- Vacationing is a multi-faceted experience with numerous benefits that go beyond just relaxation; it nourishes the mind, body, and spirit, leading to a more fulfilled life.

2214  
2215  
2216  
2217  
2218  
2219  
2220  
2221  
2222  
2223  
2224  
2225  
2226  
2227  
2228  
2229  
2230  
2231  
2232  
2233  
2234  
2235  
2236  
2237  
2238  
2239  
2240  
2241  
2242  
2243  
2244  
2245  
2246  
2247  
2248  
2249  
2250  
2251  
2252  
2253  
2254  
2255  
2256  
2257  
2258  
2259  
2260  
2261  
2262  
2263  
2264  
2265  
2266  
2267

Going on vacation offers numerous benefits that can positively impact various aspects of your life. Here are some key points to consider:

1. **Mental Health and Well-being** Vacations provide a much-needed break from routine stressors, allowing you to recharge mentally. Time away from work and daily responsibilities can reduce anxiety, improve mood, and enhance overall mental health.
2. **Increased Productivity** Studies have shown that taking time off can lead to increased productivity. A well-rested mind is more creative and efficient, which can translate to improved performance at work once you return.
3. **Quality Family Time** Vacations offer an opportunity to spend quality time with family and friends. Engaging in new activities together can strengthen relationships and create lasting memories.
4. **Exploration and Learning** Traveling exposes you to new cultures, cuisines, and experiences, broadening your perspectives and understanding of the world. This can be both educational and enriching.
5. **Physical Health Benefits** Many vacations involve physical activities such as hiking, swimming, or exploring new places, which can improve physical health. Additionally, the relaxation aspect of a vacation can lower stress levels, benefiting overall health.
6. **New Experiences and Adventures** Vacations allow you to step out of your comfort zone and try new things, fostering personal growth and a sense of adventure.
7. **Digital Detox** Being away from the usual digital environment can help you disconnect from social media and technology, leading to a healthier relationship with digital devices.
8. **Boosts Creativity** Exposure to new environments and experiences can inspire creativity and spark new ideas, which can be beneficial both personally and professionally.
9. **Opportunity for Reflection** Time away from the hustle and bustle of daily life provides a chance for introspection and reflection, allowing you to gain clarity on personal goals and aspirations.
10. **Stress Reduction** The change of scenery and relaxed pace of a vacation can significantly reduce stress, contributing to a sense of peace and well-being.
11. **Opportunity for Personal Growth** Facing new challenges and environments can build confidence and resilience, contributing to personal development.
12. **Enhanced Relationship Satisfaction** Spending uninterrupted time with loved ones can improve communication and deepen bonds, leading to more satisfying relationships.
13. **Cultural Appreciation and Tolerance** Exposure to different cultures and ways of life can foster greater appreciation and tolerance for diversity, enhancing social skills and empathy.

2268  
2269  
2270  
2271  
2272  
2273  
2274  
2275  
2276  
2277  
2278  
2279  
2280  
2281  
2282  
2283  
2284  
2285  
2286  
2287  
2288  
2289  
2290  
2291  
2292  
2293  
2294  
2295  
2296  
2297  
2298  
2299  
2300  
2301  
2302  
2303  
2304  
2305  
2306  
2307  
2308  
2309  
2310  
2311  
2312  
2313  
2314  
2315  
2316  
2317  
2318  
2319  
2320  
2321

14. **Memory Creation** Vacations are often associated with memorable experiences and stories that you can look back on fondly, creating a rich tapestry of personal history.

15. **Motivation and Inspiration** Returning from a vacation can reignite passion and motivation in various areas of life, from work to personal projects.

By considering these points, you can see how going on vacation can be a valuable investment in your overall quality of life.

Now, please respond to this question: What's the one best thing about going on a vacation?

**OLMo-2-1124-13B-Instruct:** The one best thing about going on a vacation is the profound impact it has on **mindfulness and presence**. When you step away from the familiar and immerse yourself in new experiences, environments, and cultures, you naturally become more present in the moment. This heightened state of attentiveness allows you to fully engage with your surroundings, savor the experiences, and connect more deeply with others.

Being away from the constant distractions of daily life, such as work emails, social media, and the relentless ticking of the clock, brings a sense of liberation. This freedom fosters an environment where you can truly appreciate the here and now, whether you're watching the sun set over a distant horizon, tasting unfamiliar flavors in exotic foods, or sharing a moment of pure joy with loved ones. This heightened awareness and appreciation for the present moment not only enriches the vacation experience but also carries over into your everyday life, fostering greater happiness, contentment, and a deeper connection with the world around you.

This mindfulness achieved through vacationing can lead to a more grounded and fulfilling existence, even once the vacation ends. The ability to return to life with a renewed focus on the present, an appreciation for small moments, and a more balanced perspective on life's demands is truly invaluable. This shift in consciousness is the one best thing about going on a vacation, as it has the potential to transform how you live and perceive the world.

## D MODEL PERFORMANCE DURING GAMES

2322  
2323  
2324  
2325  
2326  
2327  
2328  
2329  
2330  
2331  
2332  
2333  
2334  
2335  
2336  
2337  
2338  
2339  
2340  
2341  
2342  
2343  
2344  
2345  
2346  
2347  
2348  
2349  
2350  
2351  
2352  
2353  
2354  
2355  
2356  
2357  
2358  
2359  
2360  
2361  
2362  
2363  
2364  
2365  
2366  
2367  
2368  
2369  
2370  
2371  
2372  
2373  
2374  
2375

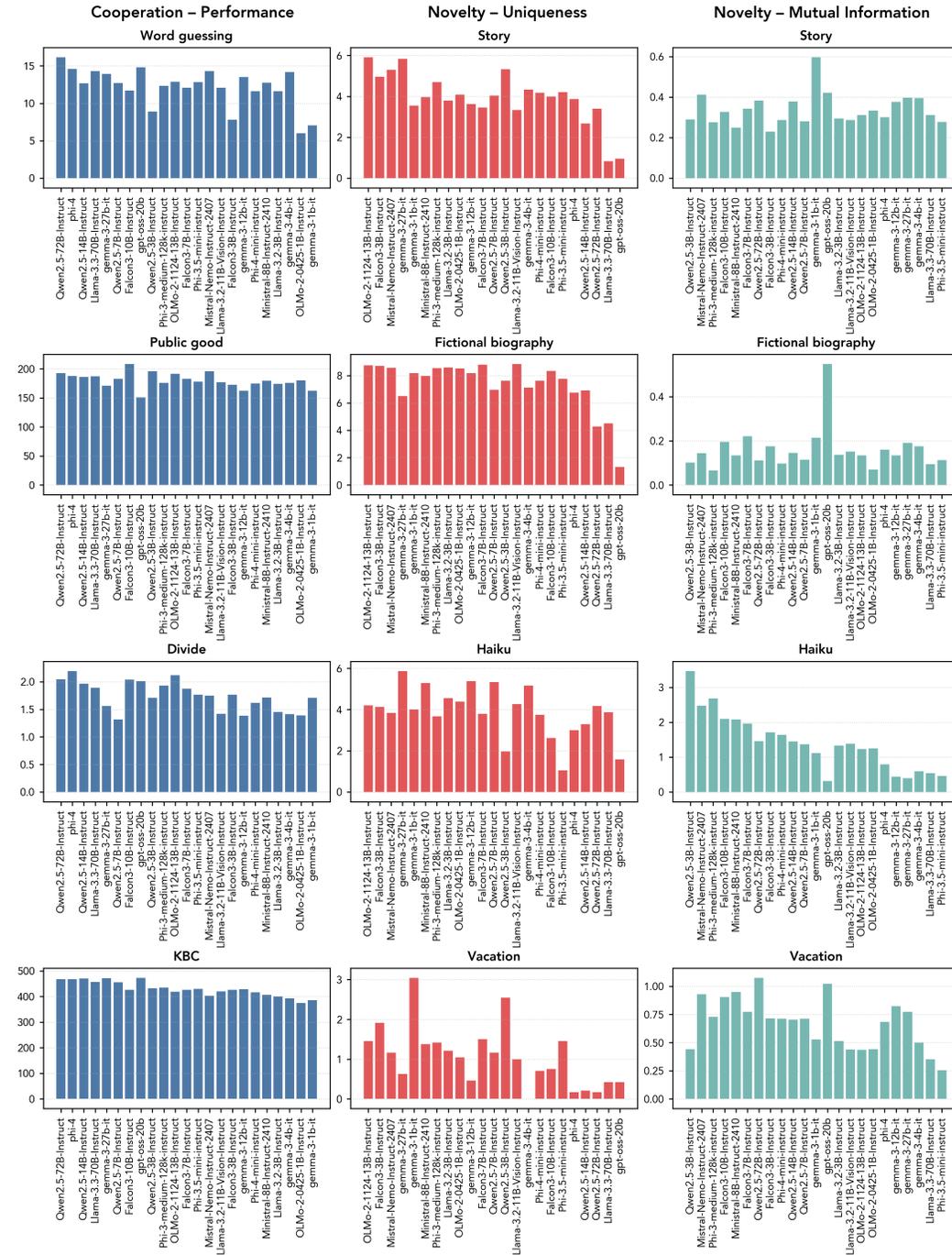


Figure 9: Average game outcomes for each model, which is calculated by averaging the model’s game outcomes across all partners (i.e., 23 models).

## E DETAILED RESULTS

### E.1 MIXED-EFFECTS REGRESSION RESULTS FOR COOPERATION

Tables 2~9 show strong positive trends across all datasets, CKA variants, and games.

#### E.1.1 WHEN USING THE GLOBAL AVERAGE

Table 2: Word guessing game. We fit a mixed-effects regression between game outcome and representational similarity. Here, the summary similarity score is calculated using a global average.  $\beta$  indicates the coefficient of similarity, and the 95% CI represents the confidence interval of  $\beta$ . **The table shows that representational similarity has a significant positive effect on the cooperative outcome.**

Dataset	Metric	Intercept	$\beta$	95% CI	$p$
WikiText	CKA (linear)	8.19	7.23	[5.45, 9.01]	$1.5 \times 10^{-15}$
	unbiased CKA (linear)	8.26	7.16	[5.40, 8.92]	$1.6 \times 10^{-15}$
	CKA (RBF)	8.26	7.33	[5.60, 9.06]	$9.0 \times 10^{-17}$
	unbiased CKA (RBF)	7.76	7.98	[6.02, 9.94]	$1.5 \times 10^{-15}$
GSM8K	CKA (linear)	10.58	6.13	[4.88, 7.38]	$6.9 \times 10^{-22}$
	unbiased CKA (linear)	10.67	6.03	[4.80, 7.26]	$9.3 \times 10^{-22}$
	CKA (RBF)	10.63	5.90	[4.68, 7.12]	$3.4 \times 10^{-21}$
	unbiased CKA (RBF)	10.55	6.27	[4.97, 7.56]	$2.4 \times 10^{-21}$
MATH	CKA (linear)	9.85	6.24	[4.94, 7.53]	$3.2 \times 10^{-21}$
	unbiased CKA (linear)	9.95	6.13	[4.86, 7.41]	$4.3 \times 10^{-21}$
	CKA (RBF)	9.71	6.25	[4.92, 7.57]	$2.5 \times 10^{-20}$
	unbiased CKA (RBF)	9.75	6.51	[5.16, 7.87]	$4.7 \times 10^{-21}$
TruthfulQA	CKA (linear)	10.03	5.96	[4.58, 7.33]	$1.8 \times 10^{-17}$
	unbiased CKA (linear)	10.14	5.90	[4.54, 7.25]	$1.3 \times 10^{-17}$
	CKA (RBF)	10.11	5.75	[4.39, 7.12]	$1.6 \times 10^{-16}$
	unbiased CKA (RBF)	9.93	6.21	[4.76, 7.66]	$4.8 \times 10^{-17}$

Table 3: Public good game. We fit a mixed-effects regression between game outcome and representational similarity. Here, the summary similarity score is calculated using a global average.  $\beta$  indicates the coefficient of similarity, and the 95% CI represents the confidence interval of  $\beta$ . **The table shows that representational similarity has a significant positive effect on the cooperative outcome.**

Dataset	Metric	Intercept	$\beta$	95% CI	$p$
WikiText	CKA (linear)	148.95	51.77	[30.82, 72.73]	$1.3 \times 10^{-6}$
	unbiased CKA (linear)	149.6	51.03	[30.44, 71.62]	$1.2 \times 10^{-6}$
	CKA (RBF)	148.3	54.36	[32.95, 75.77]	$6.5 \times 10^{-7}$
	unbiased CKA (RBF)	149.4	51.31	[29.66, 72.96]	$3.4 \times 10^{-6}$
GSM8K	CKA (linear)	162.99	52.81	[34.48, 71.14]	$1.6 \times 10^{-8}$
	unbiased CKA (linear)	163.8	52.03	[33.84, 70.21]	$2.1 \times 10^{-8}$
	CKA (RBF)	162.9	52.44	[34.45, 70.44]	$1.1 \times 10^{-8}$
	unbiased CKA (RBF)	163.1	52.82	[34.14, 71.51]	$3.0 \times 10^{-8}$
MATH	CKA (linear)	157.53	52.01	[33.91, 70.10]	$1.8 \times 10^{-8}$
	unbiased CKA (linear)	158.3	51.31	[33.42, 69.21]	$1.9 \times 10^{-8}$
	CKA (RBF)	156.0	52.85	[34.48, 71.23]	$1.7 \times 10^{-8}$
	unbiased CKA (RBF)	157.6	52.25	[33.65, 70.84]	$3.6 \times 10^{-8}$
TruthfulQA	CKA (linear)	159.93	47.62	[29.16, 66.08]	$4.3 \times 10^{-7}$
	unbiased CKA (linear)	160.8	47.06	[28.78, 65.34]	$4.5 \times 10^{-7}$
	CKA (RBF)	160.7	45.83	[27.64, 64.02]	$7.9 \times 10^{-7}$
	unbiased CKA (RBF)	160.1	47.55	[28.61, 66.49]	$8.7 \times 10^{-7}$

2430  
2431  
2432  
2433  
2434  
2435  
2436  
2437  
2438  
2439  
2440  
2441  
2442  
2443  
2444  
2445  
2446  
2447  
2448  
2449  
2450  
2451  
2452  
2453  
2454  
2455  
2456  
2457  
2458  
2459  
2460  
2461  
2462  
2463  
2464  
2465  
2466  
2467  
2468  
2469  
2470  
2471  
2472  
2473  
2474  
2475  
2476  
2477  
2478  
2479  
2480  
2481  
2482  
2483

Table 4: Divide-a-dollar game. We fit a mixed-effects regression between game outcome and representational similarity. Here, the summary similarity score is calculated using a global average.  $\beta$  indicates the coefficient of similarity, and the 95% CI represents the confidence interval of  $\beta$ . **The table shows that representational similarity has a significant positive effect on the cooperative outcome.**

Dataset	Metric	Intercept	$\beta$	95% CI	$p$
WikiText	CKA (linear)	1.47	0.44	[0.21, 0.67]	0.00014
	unbiased CKA (linear)	1.47	0.44	[0.22, 0.67]	0.00011
	CKA (RBF)	1.48	0.45	[0.23, 0.67]	$7.7 \times 10^{-5}$
	unbiased CKA (RBF)	1.46	0.46	[0.22, 0.70]	0.00017
GSM8K	CKA (linear)	1.66	0.25	[0.08, 0.43]	0.0046
	unbiased CKA (linear)	1.66	0.25	[0.08, 0.43]	0.0040
	CKA (RBF)	1.65	0.26	[0.09, 0.44]	0.0027
	unbiased CKA (RBF)	1.66	0.25	[0.07, 0.43]	0.0074
MATH	CKA (linear)	1.63	0.26	[0.08, 0.44]	0.0046
	unbiased CKA (linear)	1.63	0.27	[0.09, 0.45]	0.0036
	CKA (RBF)	1.61	0.28	[0.10, 0.47]	0.0027
	unbiased CKA (RBF)	1.63	0.26	[0.07, 0.45]	0.0066
TruthfulQA	CKA (linear)	1.64	0.23	[0.04, 0.41]	0.0185
	unbiased CKA (linear)	1.64	0.23	[0.05, 0.42]	0.0140
	CKA (RBF)	1.63	0.25	[0.06, 0.43]	0.0090
	unbiased CKA (RBF)	1.65	0.21	[0.01, 0.41]	0.0352

Table 5: KBC game. We fit a mixed-effects regression between game outcome and representational similarity. Here, the summary similarity score is calculated using a global average.  $\beta$  indicates the coefficient of similarity, and the 95% CI represents the confidence interval of  $\beta$ . **The table shows that representational similarity has a significant positive effect on the cooperative outcome.**

Dataset	Metric	Intercept	$\beta$	95% CI	$p$
WikiText	CKA (linear)	418.64	18.72	[5.75, 31.68]	0.0047
	unbiased CKA (linear)	418.8	18.53	[5.70, 31.36]	0.0046
	CKA (RBF)	419.4	17.94	[4.77, 31.10]	0.0076
	unbiased CKA (RBF)	419.3	17.77	[4.47, 31.07]	0.0088
GSM8K	CKA (linear)	423.49	19.78	[7.73, 31.84]	0.0013
	unbiased CKA (linear)	423.7	19.82	[7.89, 31.76]	0.0011
	CKA (RBF)	423.8	18.60	[6.76, 30.44]	0.0021
	unbiased CKA (RBF)	423.7	19.30	[6.99, 31.60]	0.0021
MATH	CKA (linear)	422.22	17.71	[5.90, 29.52]	0.0033
	unbiased CKA (linear)	422.4	17.68	[5.84, 29.51]	0.0034
	CKA (RBF)	422.0	17.31	[5.34, 29.27]	0.0046
	unbiased CKA (RBF)	422.5	17.15	[5.06, 29.24]	0.0054
TruthfulQA	CKA (linear)	423.78	14.46	[2.33, 26.60]	0.0195
	unbiased CKA (linear)	423.9	14.69	[2.69, 26.69]	0.0164
	CKA (RBF)	424.1	13.68	[1.75, 25.61]	0.0246
	unbiased CKA (RBF)	424.3	13.35	[0.90, 25.80]	0.0356

## E.1.2 WHEN USING THE AVERAGE OF MAXIMUM-ALIGNED SCORES

Table 6: Word guessing game. We fit a mixed-effects regression between game outcome and representational similarity. Here, the summary similarity score is calculated using the average of maximum-aligned scores.  $\beta$  indicates the coefficient of similarity, and the 95% CI represents the confidence interval of  $\beta$ . **The table shows that representational similarity has a significant positive effect on the cooperative outcome.**

Dataset	Metric	Intercept	$\beta$	95% CI	$p$
WikiText	CKA (linear)	8.37	5.58	[4.14, 7.02]	$3.5 \times 10^{-14}$
	unbiased CKA (linear)	5.78	9.10	[7.15, 11.05]	$6.3 \times 10^{-20}$
	CKA (RBF)	8.31	5.76	[4.36, 7.16]	$6.4 \times 10^{-16}$
	unbiased CKA (RBF)	8.44	5.64	[4.27, 7.01]	$7.1 \times 10^{-16}$
GSM8K	CKA (linear)	10.64	3.79	[3.00, 4.58]	$3.6 \times 10^{-21}$
	unbiased CKA (linear)	10.30	4.49	[3.62, 5.37]	$1.3 \times 10^{-23}$
	CKA (RBF)	10.50	3.92	[3.11, 4.73]	$2.5 \times 10^{-21}$
	unbiased CKA (RBF)	10.61	3.80	[3.02, 4.59]	$3.2 \times 10^{-21}$
MATH	CKA (linear)	9.96	4.32	[3.36, 5.27]	$8.1 \times 10^{-19}$
	unbiased CKA (linear)	9.27	5.51	[4.40, 6.62]	$1.6 \times 10^{-22}$
	CKA (RBF)	9.68	4.56	[3.53, 5.59]	$4.2 \times 10^{-18}$
	unbiased CKA (RBF)	9.79	4.45	[3.45, 5.46]	$5.1 \times 10^{-18}$
TruthfulQA	CKA (linear)	10.24	3.84	[2.92, 4.76]	$3.9 \times 10^{-16}$
	unbiased CKA (linear)	9.53	5.05	[3.95, 6.16]	$3.2 \times 10^{-19}$
	CKA (RBF)	10.12	3.93	[2.97, 4.89]	$1.0 \times 10^{-15}$
	unbiased CKA (RBF)	10.28	3.79	[2.87, 4.71]	$7.8 \times 10^{-16}$

Table 7: Public good game. We fit a mixed-effects regression between game outcome and representational similarity. Here, the summary similarity score is calculated using the average of maximum-aligned scores.  $\beta$  indicates the coefficient of similarity, and the 95% CI represents the confidence interval of  $\beta$ . **The table shows that representational similarity has a significant positive effect on the cooperative outcome.**

Dataset	Metric	Intercept	$\beta$	95% CI	$p$
WikiText	CKA (linear)	150.11	39.99	[20.48, 59.50]	$5.9 \times 10^{-5}$
	unbiased CKA (linear)	139.0	55.08	[30.95, 79.20]	$7.7 \times 10^{-6}$
	CKA (RBF)	149.9	40.93	[21.96, 59.90]	$2.4 \times 10^{-5}$
	unbiased CKA (RBF)	150.9	39.97	[21.39, 58.54]	$2.5 \times 10^{-5}$
GSM8K	CKA (linear)	164.33	30.77	[19.18, 42.37]	$2.0 \times 10^{-7}$
	unbiased CKA (linear)	162.1	35.51	[22.64, 48.38]	$6.4 \times 10^{-8}$
	CKA (RBF)	162.8	32.58	[20.68, 44.49]	$8.2 \times 10^{-8}$
	unbiased CKA (RBF)	163.9	31.47	[19.86, 43.08]	$1.1 \times 10^{-7}$
MATH	CKA (linear)	159.77	33.66	[19.91, 47.41]	$1.6 \times 10^{-6}$
	unbiased CKA (linear)	155.9	40.41	[24.75, 56.07]	$4.2 \times 10^{-7}$
	CKA (RBF)	156.7	37.00	[22.26, 51.74]	$8.7 \times 10^{-7}$
	unbiased CKA (RBF)	157.6	36.08	[21.68, 50.48]	$9.1 \times 10^{-7}$
TruthfulQA	CKA (linear)	161.06	31.43	[18.21, 44.65]	$3.2 \times 10^{-6}$
	unbiased CKA (linear)	156.8	38.85	[23.57, 54.13]	$6.2 \times 10^{-7}$
	CKA (RBF)	159.8	32.59	[19.07, 46.12]	$2.3 \times 10^{-6}$
	unbiased CKA (RBF)	161.5	30.82	[17.83, 43.81]	$3.3 \times 10^{-6}$

2538  
2539  
2540  
2541  
2542  
2543  
2544  
2545  
2546  
2547  
2548  
2549  
2550  
2551  
2552  
2553  
2554  
2555  
2556  
2557  
2558  
2559  
2560  
2561  
2562  
2563  
2564  
2565  
2566  
2567  
2568  
2569  
2570  
2571  
2572  
2573  
2574  
2575  
2576  
2577  
2578  
2579  
2580  
2581  
2582  
2583  
2584  
2585  
2586  
2587  
2588  
2589  
2590  
2591

Table 8: Divide-a-dollar game. We fit a mixed-effects regression between game outcome and representational similarity. Here, the summary similarity score is calculated using the average of maximum-aligned scores.  $\beta$  indicates the coefficient of similarity, and the 95% CI represents the confidence interval of  $\beta$ . **The table shows that representational similarity has a significant positive effect on the cooperative outcome.**

Dataset	Metric	Intercept	$\beta$	95% CI	$p$
WikiText	CKA (linear)	1.52	0.29	[0.11, 0.48]	0.00170
	unbiased CKA (linear)	1.46	0.38	[0.13, 0.62]	0.0025
	CKA (RBF)	1.53	0.29	[0.11, 0.47]	0.0015
	unbiased CKA (RBF)	1.53	0.28	[0.11, 0.46]	0.0015
GSM8K	CKA (linear)	1.67	0.14	[0.03, 0.24]	0.0089
	unbiased CKA (linear)	1.67	0.14	[0.03, 0.26]	0.0161
	CKA (RBF)	1.66	0.15	[0.04, 0.25]	0.0079
	unbiased CKA (RBF)	1.67	0.14	[0.04, 0.25]	0.0078
MATH	CKA (linear)	1.64	0.16	[0.03, 0.29]	0.0132
	unbiased CKA (linear)	1.63	0.17	[0.03, 0.32]	0.0204
	CKA (RBF)	1.63	0.18	[0.04, 0.31]	0.0107
	unbiased CKA (RBF)	1.63	0.17	[0.04, 0.30]	0.0105
TruthfulQA	CKA (linear)	1.64	0.16	[0.04, 0.29]	0.0082
	unbiased CKA (linear)	1.63	0.18	[0.03, 0.32]	0.0187
	CKA (RBF)	1.63	0.18	[0.06, 0.31]	0.0045
	unbiased CKA (RBF)	1.64	0.17	[0.05, 0.29]	0.0050

Table 9: KBC game. We fit a mixed-effects regression between game outcome and representational similarity. Here, the summary similarity score is calculated using the average of maximum-aligned scores.  $\beta$  indicates the coefficient of similarity, and the 95% CI represents the confidence interval of  $\beta$ . **The table shows that representational similarity has a significant positive effect on the cooperative outcome.**

Dataset	Metric	Intercept	$\beta$	95% CI	$p$
WikiText	CKA (linear)	413.25	22.15	[9.64, 34.65]	0.00052
	unbiased CKA (linear)	411.2	25.12	[9.80, 40.44]	0.00131
	CKA (RBF)	414.1	21.43	[9.13, 33.72]	0.00064
	unbiased CKA (RBF)	414.5	21.01	[8.96, 33.07]	0.00063
GSM8K	CKA (linear)	422.86	13.71	[6.19, 21.24]	0.00036
	unbiased CKA (linear)	422.4	14.82	[6.46, 23.18]	0.00051
	CKA (RBF)	422.8	13.44	[5.70, 21.17]	0.00066
	unbiased CKA (RBF)	423.0	13.35	[5.81, 20.89]	0.00052
MATH	CKA (linear)	420.55	15.44	[6.50, 24.39]	0.00072
	unbiased CKA (linear)	419.6	17.18	[6.99, 27.37]	0.00095
	CKA (RBF)	419.8	15.88	[6.29, 25.47]	0.00117
	unbiased CKA (RBF)	420.1	15.75	[6.38, 25.12]	0.00098
TruthfulQA	CKA (linear)	421.48	13.88	[5.30, 22.46]	0.00152
	unbiased CKA (linear)	420.9	14.97	[5.02, 24.91]	0.00318
	CKA (RBF)	421.4	13.67	[4.88, 22.46]	0.00230
	unbiased CKA (RBF)	421.8	13.49	[5.05, 21.92]	0.00172

## E.2 MIXED-EFFECTS REGRESSION RESULTS FOR UNIQUENESS

Tables 10~17 consistently show negative trends across all datasets, CKA variants, and games.

## E.2.1 WHEN USING THE GLOBAL AVERAGE

Table 10: Story writing task. We fit a mixed-effects regression between response uniqueness and representational similarity. Here, the summary similarity score is calculated using a global average.  $\beta$  indicates the coefficient of similarity, and the 95% CI represents the confidence interval of  $\beta$ . **The table shows that representational similarity has a consistent negative effect on response uniqueness.**

Dataset	Metric	Intercept	$\beta$	95% CI	$p$
WikiText	CKA (linear)	4.12	-0.31	[-1.56, 0.95]	0.633
	unbiased CKA (linear)	4.12	-0.30	[-1.55, 0.94]	0.630
	CKA (RBF)	4.14	-0.34	[-1.63, 0.96]	0.610
	unbiased CKA (RBF)	4.14	-0.33	[-1.60, 0.93]	0.608
GSM8K	CKA (linear)	4.03	-0.27	[-1.42, 0.88]	0.643
	unbiased CKA (linear)	4.02	-0.26	[-1.40, 0.88]	0.659
	CKA (RBF)	4.04	-0.29	[-1.44, 0.86]	0.626
	unbiased CKA (RBF)	4.03	-0.26	[-1.39, 0.87]	0.655
MATH	CKA (linear)	4.10	-0.36	[-1.48, 0.77]	0.536
	unbiased CKA (linear)	4.09	-0.34	[-1.45, 0.77]	0.548
	CKA (RBF)	4.12	-0.37	[-1.54, 0.79]	0.532
	unbiased CKA (RBF)	4.10	-0.35	[-1.49, 0.80]	0.553
TruthfulQA	CKA (linear)	4.15	-0.49	[-1.65, 0.66]	0.403
	unbiased CKA (linear)	4.13	-0.46	[-1.61, 0.68]	0.428
	CKA (RBF)	4.14	-0.45	[-1.61, 0.72]	0.453
	unbiased CKA (RBF)	4.11	-0.40	[-1.54, 0.74]	0.490

Table 11: Fictional biography generation task. We fit a mixed-effects regression between response uniqueness and representational similarity. Here, the summary similarity score is calculated using a global average.  $\beta$  indicates the coefficient of similarity, and the 95% CI represents the confidence interval of  $\beta$ . **The table shows that representational similarity has a significant negative effect on response uniqueness.**

Dataset	Metric	Intercept	$\beta$	95% CI	$p$
WikiText	CKA (linear)	8.14	-1.28	[-2.47, -0.10]	0.034
	unbiased CKA (linear)	8.13	-1.28	[-2.45, -0.11]	0.0327
	CKA (RBF)	8.21	-1.42	[-2.59, -0.24]	0.0181
	unbiased CKA (RBF)	8.19	-1.40	[-2.55, -0.25]	0.0169
GSM8K	CKA (linear)	7.94	-1.75	[-2.63, -0.87]	$9.3 \times 10^{-5}$
	unbiased CKA (linear)	7.92	-1.76	[-2.63, -0.90]	$6.7 \times 10^{-5}$
	CKA (RBF)	7.95	-1.67	[-2.55, -0.79]	$1.89 \times 10^{-4}$
	unbiased CKA (RBF)	7.92	-1.69	[-2.55, -0.83]	$1.16 \times 10^{-4}$
MATH	CKA (linear)	8.07	-1.60	[-2.52, -0.69]	$5.7 \times 10^{-4}$
	unbiased CKA (linear)	8.05	-1.61	[-2.51, -0.71]	$4.60 \times 10^{-4}$
	CKA (RBF)	8.12	-1.61	[-2.56, -0.65]	$9.47 \times 10^{-4}$
	unbiased CKA (RBF)	8.10	-1.61	[-2.54, -0.68]	$7.17 \times 10^{-4}$
TruthfulQA	CKA (linear)	7.92	-1.29	[-2.24, -0.35]	0.0073
	unbiased CKA (linear)	7.91	-1.33	[-2.27, -0.40]	0.00506
	CKA (RBF)	7.91	-1.21	[-2.17, -0.25]	0.0134
	unbiased CKA (RBF)	7.90	-1.26	[-2.19, -0.32]	0.00824

2646  
2647  
2648  
2649  
2650  
2651  
2652  
2653  
2654  
2655  
2656  
2657  
2658  
2659  
2660  
2661  
2662  
2663  
2664  
2665  
2666  
2667  
2668  
2669  
2670  
2671  
2672  
2673  
2674  
2675  
2676  
2677  
2678  
2679  
2680  
2681  
2682  
2683  
2684  
2685  
2686  
2687  
2688  
2689  
2690  
2691  
2692  
2693  
2694  
2695  
2696  
2697  
2698  
2699

Table 12: Haiku composition task. We fit a mixed-effects regression between response uniqueness and representational similarity. Here, the summary similarity score is calculated using a global average.  $\beta$  indicates the coefficient of similarity, and the 95% CI represents the confidence interval of  $\beta$ . **The table shows that representational similarity has a significant negative effect on response uniqueness.**

Dataset	Metric	Intercept	$\beta$	95% CI	$p$
WikiText	CKA (linear)	5.95	-3.42	[-4.80, -2.05]	$1.1 \times 10^{-6}$
	unbiased CKA (linear)	5.93	-3.42	[-4.77, -2.07]	$6.71 \times 10^{-7}$
	CKA (RBF)	5.74	-3.13	[-4.56, -1.70]	$1.71 \times 10^{-5}$
	unbiased CKA (RBF)	5.70	-3.10	[-4.49, -1.70]	$1.39 \times 10^{-5}$
GSM8K	CKA (linear)	4.67	-2.43	[-3.70, -1.16]	$1.7 \times 10^{-4}$
	unbiased CKA (linear)	4.64	-2.42	[-3.68, -1.17]	0.000158
	CKA (RBF)	4.65	-2.22	[-3.49, -0.96]	0.000589
	unbiased CKA (RBF)	4.61	-2.22	[-3.46, -0.97]	0.000477
MATH	CKA (linear)	5.24	-3.12	[-4.36, -1.89]	$7.2 \times 10^{-7}$
	unbiased CKA (linear)	5.22	-3.13	[-4.35, -1.92]	$4.48 \times 10^{-7}$
	CKA (RBF)	5.39	-3.19	[-4.45, -1.93]	$6.88 \times 10^{-7}$
	unbiased CKA (RBF)	5.33	-3.17	[-4.40, -1.94]	$4.60 \times 10^{-7}$
TruthfulQA	CKA (linear)	4.98	-2.58	[-3.84, -1.33]	$5.2 \times 10^{-5}$
	unbiased CKA (linear)	4.94	-2.58	[-3.83, -1.34]	$4.75 \times 10^{-5}$
	CKA (RBF)	4.84	-2.12	[-3.43, -0.81]	0.00149
	unbiased CKA (RBF)	4.95	-2.19	[-3.83, -0.91]	0.000794

Table 13: Vacation brainstorming task. We fit a mixed-effects regression between response uniqueness and representational similarity. Here, the summary similarity score is calculated using a global average.  $\beta$  indicates the coefficient of similarity, and the 95% CI represents the confidence interval of  $\beta$ . **The table shows that representational similarity has a significant negative effect on response uniqueness.**

Dataset	Metric	Intercept	$\beta$	95% CI	$p$
WikiText	CKA (linear)	1.66	-1.01	[-1.58, -0.45]	0.00045
	unbiased CKA (linear)	1.65	-1.00	[-1.56, -0.44]	0.000480
	CKA (RBF)	1.61	-0.94	[-1.55, -0.32]	0.00287
	unbiased CKA (RBF)	1.59	-0.91	[-1.51, -0.31]	0.00308
GSM8K	CKA (linear)	1.26	-0.65	[-1.32, 0.01]	0.053
	unbiased CKA (linear)	1.25	-0.64	[-1.29, 0.02]	0.0584
	CKA (RBF)	1.26	-0.60	[-1.26, 0.05]	0.0721
	unbiased CKA (RBF)	1.24	-0.57	[-1.22, 0.07]	0.0824
MATH	CKA (linear)	1.45	-0.91	[-1.50, -0.31]	0.0028
	unbiased CKA (linear)	1.43	-0.89	[-1.48, -0.30]	0.00298
	CKA (RBF)	1.45	-0.87	[-1.47, -0.27]	0.00468
	unbiased CKA (RBF)	1.45	-0.87	[-1.47, -0.27]	0.00468
TruthfulQA	CKA (linear)	1.35	-0.71	[-1.34, -0.08]	0.0273
	unbiased CKA (linear)	1.34	-0.69	[-1.31, -0.07]	0.0304
	CKA (RBF)	1.34	-0.65	[-1.27, -0.02]	0.0438
	unbiased CKA (RBF)	1.32	-0.62	[-1.24, -0.01]	0.0480

## E.2.2 WHEN USING THE AVERAGE OF MAXIMUM-ALIGNED SCORES

Table 14: Story writing task. We fit a mixed-effects regression between response uniqueness and representational similarity. Here, the summary similarity score is calculated using the average of maximum-aligned scores.  $\beta$  indicates the coefficient of similarity, and the 95% CI represents the confidence interval of  $\beta$ . **The table shows that representational similarity has a consistent negative effect on response uniqueness.**

Dataset	Metric	Intercept	$\beta$	95% CI	$p$
WikiText	CKA (linear)	4.31	-0.49	[-1.69, 0.71]	0.422
	unbiased CKA (linear)	4.31	-0.50	[-1.68, 0.69]	0.412
	CKA (RBF)	4.31	-0.51	[-1.66, 0.65]	0.391
	unbiased CKA (RBF)	4.31	-0.50	[-1.68, 0.67]	0.402
GSM8K	CKA (linear)	4.02	-0.16	[-0.88, 0.56]	0.656
	unbiased CKA (linear)	4.02	-0.16	[-0.87, 0.55]	0.662
	CKA (RBF)	4.03	-0.17	[-0.89, 0.55]	0.641
	unbiased CKA (RBF)	4.04	-0.18	[-0.92, 0.56]	0.632
MATH	CKA (linear)	4.11	-0.29	[-1.14, 0.57]	0.512
	unbiased CKA (linear)	4.11	-0.28	[-1.13, 0.56]	0.514
	CKA (RBF)	4.13	-0.31	[-1.20, 0.59]	0.502
	unbiased CKA (RBF)	4.14	-0.32	[-1.23, 0.60]	0.500
TruthfulQA	CKA (linear)	4.06	-0.19	[-1.01, 0.63]	0.644
	unbiased CKA (linear)	4.05	-0.19	[-0.99, 0.61]	0.641
	CKA (RBF)	4.05	-0.19	[-1.00, 0.62]	0.646
	unbiased CKA (RBF)	4.06	-0.19	[-1.03, 0.65]	0.654

Table 15: Fictional biography generation task. We fit a mixed-effects regression between response uniqueness and representational similarity. Here, the summary similarity score is calculated using the average of maximum-aligned scores.  $\beta$  indicates the coefficient of similarity, and the 95% CI represents the confidence interval of  $\beta$ . **The table shows that representational similarity has a significant negative effect on response uniqueness.**

Dataset	Metric	Intercept	$\beta$	95% CI	$p$
WikiText	CKA (linear)	8.53	-1.55	[-2.47, -0.63]	0.00094
	unbiased CKA (linear)	8.52	-1.54	[-2.45, -0.63]	0.001
	CKA (RBF)	8.50	-1.55	[-2.43, -0.68]	0.001
	unbiased CKA (RBF)	8.53	-1.58	[-2.47, -0.68]	0.001
GSM8K	CKA (linear)	7.96	-1.15	[-1.67, -0.62]	$1.8 \times 10^{-5}$
	unbiased CKA (linear)	7.95	-1.14	[-1.66, -0.63]	$1.0 \times 10^{-5}$
	CKA (RBF)	7.96	-1.15	[-1.67, -0.62]	$2.0 \times 10^{-5}$
	unbiased CKA (RBF)	7.98	-1.16	[-1.70, -0.62]	$3.0 \times 10^{-5}$
MATH	CKA (linear)	8.17	-1.31	[-1.94, -0.69]	$4.3 \times 10^{-5}$
	unbiased CKA (linear)	8.15	-1.31	[-1.93, -0.69]	$4.0 \times 10^{-5}$
	CKA (RBF)	8.23	-1.37	[-2.03, -0.71]	$5.0 \times 10^{-5}$
	unbiased CKA (RBF)	8.25	-1.39	[-2.07, -0.71]	$6.0 \times 10^{-5}$
TruthfulQA	CKA (linear)	8.08	-1.18	[-1.79, -0.56]	0.00016
	unbiased CKA (linear)	8.06	-1.16	[-1.76, -0.56]	$1.0 \times 10^{-4}$
	CKA (RBF)	8.04	-1.12	[-1.72, -0.51]	$3.0 \times 10^{-4}$
	unbiased CKA (RBF)	8.08	-1.14	[-1.78, -0.51]	$4.0 \times 10^{-4}$

2754  
2755  
2756  
2757  
2758  
2759  
2760  
2761  
2762  
2763  
2764  
2765  
2766  
2767  
2768  
2769  
2770  
2771  
2772  
2773  
2774  
2775  
2776  
2777  
2778  
2779  
2780  
2781  
2782  
2783  
2784  
2785  
2786  
2787  
2788  
2789  
2790  
2791  
2792  
2793  
2794  
2795  
2796  
2797  
2798  
2799  
2800  
2801  
2802  
2803  
2804  
2805  
2806  
2807

Table 16: Haiku composition task. We fit a mixed-effects regression between response uniqueness and representational similarity. Here, the summary similarity score is calculated using the average of maximum-aligned scores.  $\beta$  indicates the coefficient of similarity, and the 95% CI represents the confidence interval of  $\beta$ . **The table shows that representational similarity has a significant negative effect on response uniqueness.**

Dataset	Metric	Intercept	$\beta$	95% CI	$p$
WikiText	CKA (linear)	5.86	-2.64	[-3.95, -1.32]	$8.9 \times 10^{-5}$
	unbiased CKA (linear)	5.84	-2.61	[-3.92, -1.31]	$8.75 \times 10^{-5}$
	CKA (RBF)	5.62	-2.37	[-3.63, -1.11]	$2.33 \times 10^{-4}$
	unbiased CKA (RBF)	5.65	-2.39	[-3.68, -1.11]	$2.67 \times 10^{-4}$
GSM8K	CKA (linear)	4.62	-1.44	[-2.22, -0.66]	0.00028
	unbiased CKA (linear)	4.60	-1.43	[-2.19, -0.66]	$2.57 \times 10^{-4}$
	CKA (RBF)	4.59	-1.38	[-2.16, -0.60]	$5.19 \times 10^{-4}$
	unbiased CKA (RBF)	4.62	-1.40	[-2.20, -0.60]	$6.01 \times 10^{-4}$
MATH	CKA (linear)	5.04	-1.92	[-2.84, -1.00]	$4.9 \times 10^{-5}$
	unbiased CKA (linear)	5.03	-1.91	[-2.83, -1.00]	$4.09 \times 10^{-5}$
	CKA (RBF)	5.11	-1.97	[-2.94, -1.00]	$6.83 \times 10^{-5}$
	unbiased CKA (RBF)	5.14	-1.98	[-2.98, -0.99]	$9.19 \times 10^{-5}$
TruthfulQA	CKA (linear)	4.88	-1.65	[-2.54, -0.76]	0.00030
	unbiased CKA (linear)	4.86	-1.64	[-2.51, -0.76]	$2.43 \times 10^{-4}$
	CKA (RBF)	4.78	-1.49	[-2.37, -0.61]	$8.84 \times 10^{-4}$
	unbiased CKA (RBF)	4.82	-1.51	[-2.42, -0.59]	0.0013

Table 17: Vacation brainstorming task. We fit a mixed-effects regression between response uniqueness and representational similarity. Here, the summary similarity score is calculated using the average of maximum-aligned scores.  $\beta$  indicates the coefficient of similarity, and the 95% CI represents the confidence interval of  $\beta$ . **The table shows that representational similarity has a significant negative effect on response uniqueness.**

Dataset	Metric	Intercept	$\beta$	95% CI	$p$
WikiText	CKA (linear)	1.79	-0.97	[-1.65, -0.30]	0.0046
	unbiased CKA (linear)	1.78	-0.97	[-1.63, -0.30]	0.005
	CKA (RBF)	1.69	-0.87	[-1.52, -0.21]	0.009
	unbiased CKA (RBF)	1.70	-0.88	[-1.55, -0.22]	0.009
GSM8K	CKA (linear)	1.25	-0.38	[-0.80, 0.05]	0.081
	unbiased CKA (linear)	1.24	-0.37	[-0.80, 0.05]	0.083
	CKA (RBF)	1.25	-0.37	[-0.80, 0.05]	0.083
	unbiased CKA (RBF)	1.26	-0.39	[-0.82, 0.05]	0.080
MATH	CKA (linear)	1.44	-0.64	[-1.14, -0.14]	0.012
	unbiased CKA (linear)	1.44	-0.64	[-1.13, -0.14]	0.011
	CKA (RBF)	1.48	-0.68	[-1.20, -0.16]	0.011
	unbiased CKA (RBF)	1.49	-0.69	[-1.22, -0.16]	0.011
TruthfulQA	CKA (linear)	1.36	-0.50	[-0.98, -0.03]	0.039
	unbiased CKA (linear)	1.35	-0.50	[-0.96, -0.03]	0.036
	CKA (RBF)	1.36	-0.51	[-0.97, -0.04]	0.032
	unbiased CKA (RBF)	1.37	-0.52	[-1.00, -0.03]	0.036

E.3 MIXED-EFFECTS REGRESSION  
 RESULTS FOR MUTUAL INFORMATION CALCULATED WITH LLAMA-3.1-8B-INSTRUCT

Tables 18~25 show significant positive trends across all datasets, CKA variants, and games.

E.3.1 WHEN USING THE GLOBAL AVERAGE

Table 18: Story writing task. We fit a mixed-effects regression between mutual information and representational similarity. Here, the summary similarity score is calculated using a global average.  $\beta$  indicates the coefficient of similarity, and the 95% CI represents the confidence interval of  $\beta$ . **The table shows that representational similarity has a significant positive effect on mutual information.**

Dataset	Metric	Intercept	$\beta$	95% CI	$p$
WikiText	CKA (linear)	0.303	0.053	[0.020, 0.087]	0.00194
	unbiased CKA (linear)	0.304	0.053	[0.019, 0.086]	0.00193
	CKA (RBF)	0.306	0.050	[0.016, 0.084]	0.00379
	unbiased CKA (RBF)	0.307	0.049	[0.016, 0.082]	0.00372
GSM8K	CKA (linear)	0.317	0.057	[0.027, 0.087]	0.00021
	unbiased CKA (linear)	0.318	0.057	[0.027, 0.086]	0.00020
	CKA (RBF)	0.316	0.055	[0.025, 0.085]	0.00035
	unbiased CKA (RBF)	0.318	0.054	[0.025, 0.084]	0.00032
MATH	CKA (linear)	0.313	0.052	[0.022, 0.082]	0.00058
	unbiased CKA (linear)	0.313	0.052	[0.022, 0.081]	0.00055
	CKA (RBF)	0.310	0.054	[0.024, 0.085]	0.00051
	unbiased CKA (RBF)	0.311	0.053	[0.023, 0.083]	0.00047
TruthfulQA	CKA (linear)	0.312	0.056	[0.026, 0.086]	0.00026
	unbiased CKA (linear)	0.312	0.057	[0.027, 0.086]	0.00019
	CKA (RBF)	0.312	0.053	[0.023, 0.083]	0.00062
	unbiased CKA (RBF)	0.313	0.053	[0.024, 0.083]	0.00040

Table 19: Fictional biography generation task. We fit a mixed-effects regression between mutual information and representational similarity. Here, the summary similarity score is calculated using a global average.  $\beta$  indicates the coefficient of similarity, and the 95% CI represents the confidence interval of  $\beta$ . **The table shows that representational similarity has a significant positive effect on mutual information.**

Dataset	Metric	Intercept	$\beta$	95% CI	$p$
WikiText	CKA (linear)	0.083	0.121	[0.091, 0.150]	$5.07 \times 10^{-16}$
	unbiased CKA (linear)	0.084	0.120	[0.091, 0.149]	$4.71 \times 10^{-16}$
	CKA (RBF)	0.084	0.119	[0.090, 0.149]	$2.31 \times 10^{-15}$
	unbiased CKA (RBF)	0.087	0.117	[0.088, 0.146]	$1.96 \times 10^{-15}$
GSM8K	CKA (linear)	0.123	0.099	[0.076, 0.122]	$2.17 \times 10^{-17}$
	unbiased CKA (linear)	0.124	0.098	[0.075, 0.121]	$2.08 \times 10^{-17}$
	CKA (RBF)	0.121	0.098	[0.075, 0.121]	$6.39 \times 10^{-17}$
	unbiased CKA (RBF)	0.124	0.096	[0.074, 0.119]	$5.27 \times 10^{-17}$
MATH	CKA (linear)	0.109	0.107	[0.084, 0.131]	$1.39 \times 10^{-19}$
	unbiased CKA (linear)	0.110	0.106	[0.083, 0.129]	$1.13 \times 10^{-19}$
	CKA (RBF)	0.103	0.111	[0.086, 0.135]	$5.17 \times 10^{-19}$
	unbiased CKA (RBF)	0.106	0.108	[0.085, 0.132]	$3.86 \times 10^{-19}$
TruthfulQA	CKA (linear)	0.117	0.091	[0.067, 0.115]	$2.38 \times 10^{-13}$
	unbiased CKA (linear)	0.118	0.091	[0.067, 0.115]	$1.03 \times 10^{-13}$
	CKA (RBF)	0.115	0.088	[0.064, 0.113]	$2.44 \times 10^{-12}$
	unbiased CKA (RBF)	0.118	0.088	[0.064, 0.112]	$6.81 \times 10^{-13}$

2862  
2863  
2864  
2865  
2866  
2867  
2868  
2869  
2870  
2871  
2872  
2873  
2874  
2875  
2876  
2877  
2878  
2879  
2880  
2881  
2882  
2883  
2884  
2885  
2886  
2887  
2888  
2889  
2890  
2891  
2892  
2893  
2894  
2895  
2896  
2897  
2898  
2899  
2900  
2901  
2902  
2903  
2904  
2905  
2906  
2907  
2908  
2909  
2910  
2911  
2912  
2913  
2914  
2915

Table 20: Haiku composition task. We fit a mixed-effects regression between mutual information and representational similarity. Here, the summary similarity score is calculated using a global average.  $\beta$  indicates the coefficient of similarity, and the 95% CI represents the confidence interval of  $\beta$ . **The table shows that representational similarity has a significant positive effect on mutual information.**

Dataset	Metric	Intercept	$\beta$	95% CI	$p$
WikiText	CKA (linear)	0.594	1.310	[1.03, 1.59]	$1.26 \times 10^{-20}$
	unbiased CKA (linear)	0.603	1.302	[1.029, 1.575]	$9.20 \times 10^{-21}$
	CKA (RBF)	0.672	1.194	[0.919, 1.470]	$2.03 \times 10^{-17}$
	unbiased CKA (RBF)	0.693	1.177	[0.907, 1.447]	$1.28 \times 10^{-17}$
GSM8K	CKA (linear)	1.161	0.688	[0.48, 0.89]	$6.18 \times 10^{-11}$
	unbiased CKA (linear)	1.171	0.679	[0.475, 0.882]	$6.20 \times 10^{-11}$
	CKA (RBF)	1.152	0.670	[0.463, 0.877]	$2.15 \times 10^{-10}$
	unbiased CKA (RBF)	1.168	0.656	[0.454, 0.858]	$1.88 \times 10^{-10}$
MATH	CKA (linear)	1.017	0.845	[0.63, 1.06]	$8.56 \times 10^{-15}$
	unbiased CKA (linear)	1.025	0.842	[0.632, 1.053]	$4.48 \times 10^{-15}$
	CKA (RBF)	0.970	0.878	[0.654, 1.101]	$1.33 \times 10^{-14}$
	unbiased CKA (RBF)	0.986	0.871	[0.653, 1.089]	$4.73 \times 10^{-15}$
TruthfulQA	CKA (linear)	1.046	0.794	[0.57, 1.02]	$2.62 \times 10^{-12}$
	unbiased CKA (linear)	1.054	0.804	[0.585, 1.024]	$6.80 \times 10^{-13}$
	CKA (RBF)	1.039	0.766	[0.539, 0.992]	$3.51 \times 10^{-11}$
	unbiased CKA (RBF)	1.053	0.777	[0.557, 0.997]	$4.51 \times 10^{-12}$

Table 21: Vacation brainstorming task. We fit a mixed-effects regression between mutual information and representational similarity. Here, the summary similarity score is calculated using a global average.  $\beta$  indicates the coefficient of similarity, and the 95% CI represents the confidence interval of  $\beta$ . **The table shows that representational similarity has a significant positive effect on mutual information.**

Dataset	Metric	Intercept	$\beta$	95% CI	$p$
WikiText	CKA (linear)	0.587	0.132	[0.071, 0.194]	$2.40 \times 10^{-5}$
	unbiased CKA (linear)	0.589	0.131	[0.070, 0.192]	$2.43 \times 10^{-5}$
	CKA (RBF)	0.592	0.127	[0.061, 0.192]	0.00015
	unbiased CKA (RBF)	0.594	0.125	[0.060, 0.189]	0.00014
GSM8K	CKA (linear)	0.621	0.141	[0.078, 0.203]	$9.48 \times 10^{-6}$
	unbiased CKA (linear)	0.623	0.140	[0.079, 0.202]	$8.27 \times 10^{-6}$
	CKA (RBF)	0.619	0.140	[0.077, 0.202]	$1.09 \times 10^{-5}$
	unbiased CKA (RBF)	0.621	0.139	[0.078, 0.200]	$8.45 \times 10^{-6}$
MATH	CKA (linear)	0.602	0.149	[0.090, 0.208]	$7.62 \times 10^{-7}$
	unbiased CKA (linear)	0.604	0.148	[0.090, 0.206]	$6.08 \times 10^{-7}$
	CKA (RBF)	0.594	0.154	[0.093, 0.215]	$7.74 \times 10^{-7}$
	unbiased CKA (RBF)	0.597	0.153	[0.093, 0.212]	$5.49 \times 10^{-7}$
TruthfulQA	CKA (linear)	0.619	0.113	[0.051, 0.175]	0.00034
	unbiased CKA (linear)	0.621	0.112	[0.051, 0.174]	0.00032
	CKA (RBF)	0.618	0.108	[0.046, 0.170]	0.00068
	unbiased CKA (RBF)	0.622	0.107	[0.046, 0.167]	0.00060

## E.3.2 WHEN USING THE AVERAGE OF MAXIMUM-ALIGNED SCORES

Table 22: Story writing task. We fit a mixed-effects regression between mutual information and representational similarity. Here, the summary similarity score is calculated using the average of maximum-aligned scores.  $\beta$  indicates the coefficient of similarity, and the 95% CI represents the confidence interval of  $\beta$ . **The table shows that representational similarity has a significant positive effect on mutual information.**

Dataset	Metric	Intercept	$\beta$	95% CI	$p$
WikiText	CKA (linear)	0.272	0.085	[0.053, 0.117]	$2.64 \times 10^{-7}$
	unbiased CKA (linear)	0.272	0.084	[0.052, 0.116]	$2.62 \times 10^{-7}$
	CKA (RBF)	0.276	0.081	[0.049, 0.112]	$5.54 \times 10^{-7}$
	unbiased CKA (RBF)	0.277	0.079	[0.048, 0.110]	$5.23 \times 10^{-7}$
GSM8K	CKA (linear)	0.312	0.045	[0.026, 0.065]	$4.39 \times 10^{-6}$
	unbiased CKA (linear)	0.313	0.045	[0.026, 0.064]	$4.35 \times 10^{-6}$
	CKA (RBF)	0.311	0.046	[0.026, 0.066]	$5.92 \times 10^{-6}$
	unbiased CKA (RBF)	0.312	0.045	[0.025, 0.064]	$5.76 \times 10^{-6}$
MATH	CKA (linear)	0.301	0.057	[0.034, 0.080]	$9.67 \times 10^{-7}$
	unbiased CKA (linear)	0.302	0.056	[0.034, 0.079]	$9.72 \times 10^{-7}$
	CKA (RBF)	0.297	0.061	[0.036, 0.085]	$1.09 \times 10^{-6}$
	unbiased CKA (RBF)	0.298	0.059	[0.036, 0.083]	$1.09 \times 10^{-6}$
TruthfulQA	CKA (linear)	0.302	0.056	[0.034, 0.078]	$5.79 \times 10^{-7}$
	unbiased CKA (linear)	0.303	0.055	[0.033, 0.076]	$4.89 \times 10^{-7}$
	CKA (RBF)	0.301	0.055	[0.033, 0.078]	$1.35 \times 10^{-6}$
	unbiased CKA (RBF)	0.303	0.054	[0.032, 0.075]	$9.85 \times 10^{-7}$

Table 23: Fictional biography generation task. We fit a mixed-effects regression between mutual information and representational similarity. Here, the summary similarity score is calculated using the average of maximum-aligned scores.  $\beta$  indicates the coefficient of similarity, and the 95% CI represents the confidence interval of  $\beta$ . **The table shows that representational similarity has a significant positive effect on mutual information.**

Dataset	Metric	Intercept	$\beta$	95% CI	$p$
WikiText	CKA (linear)	0.089	0.088	[0.064, 0.112]	$6.16 \times 10^{-13}$
	unbiased CKA (linear)	0.090	0.087	[0.063, 0.111]	$7.49 \times 10^{-13}$
	CKA (RBF)	0.091	0.088	[0.064, 0.111]	$1.86 \times 10^{-13}$
	unbiased CKA (RBF)	0.093	0.086	[0.063, 0.109]	$2.40 \times 10^{-13}$
GSM8K	CKA (linear)	0.124	0.061	[0.047, 0.075]	$2.48 \times 10^{-18}$
	unbiased CKA (linear)	0.125	0.060	[0.047, 0.074]	$3.41 \times 10^{-18}$
	CKA (RBF)	0.121	0.064	[0.050, 0.078]	$1.05 \times 10^{-18}$
	unbiased CKA (RBF)	0.123	0.062	[0.048, 0.075]	$1.70 \times 10^{-18}$
MATH	CKA (linear)	0.112	0.071	[0.055, 0.088]	$2.68 \times 10^{-17}$
	unbiased CKA (linear)	0.113	0.070	[0.054, 0.086]	$4.29 \times 10^{-17}$
	CKA (RBF)	0.107	0.077	[0.059, 0.094]	$2.40 \times 10^{-17}$
	unbiased CKA (RBF)	0.109	0.074	[0.057, 0.091]	$5.19 \times 10^{-17}$
TruthfulQA	CKA (linear)	0.116	0.065	[0.049, 0.081]	$2.18 \times 10^{-15}$
	unbiased CKA (linear)	0.118	0.063	[0.047, 0.078]	$2.94 \times 10^{-15}$
	CKA (RBF)	0.114	0.065	[0.049, 0.082]	$6.08 \times 10^{-15}$
	unbiased CKA (RBF)	0.118	0.063	[0.047, 0.078]	$7.99 \times 10^{-15}$

2970  
2971  
2972  
2973  
2974  
2975  
2976  
2977  
2978  
2979  
2980  
2981  
2982  
2983  
2984  
2985  
2986  
2987  
2988  
2989  
2990  
2991  
2992  
2993  
2994  
2995  
2996  
2997  
2998  
2999  
3000  
3001  
3002  
3003  
3004  
3005  
3006  
3007  
3008  
3009  
3010  
3011  
3012  
3013  
3014  
3015  
3016  
3017  
3018  
3019  
3020  
3021  
3022  
3023

Table 24: Haiku composition task. We fit a mixed-effects regression between mutual information and representational similarity. Here, the summary similarity score is calculated using the average of maximum-aligned scores.  $\beta$  indicates the coefficient of similarity, and the 95% CI represents the confidence interval of  $\beta$ . **The table shows that representational similarity has a significant positive effect on mutual information.**

Dataset	Metric	Intercept	$\beta$	95% CI	$p$
WikiText	CKA (linear)	0.652	0.973	[0.760, 1.187]	$4.43 \times 10^{-19}$
	unbiased CKA (linear)	0.659	0.966	[0.754, 1.178]	$3.84 \times 10^{-19}$
	CKA (RBF)	0.714	0.907	[0.699, 1.115]	$1.29 \times 10^{-17}$
	unbiased CKA (RBF)	0.729	0.893	[0.689, 1.097]	$9.53 \times 10^{-18}$
GSM8K	CKA (linear)	1.123	0.507	[0.385, 0.629]	$3.05 \times 10^{-16}$
	unbiased CKA (linear)	1.131	0.501	[0.381, 0.620]	$2.77 \times 10^{-16}$
	CKA (RBF)	1.111	0.513	[0.387, 0.638]	$1.05 \times 10^{-15}$
	unbiased CKA (RBF)	1.125	0.500	[0.378, 0.622]	$9.96 \times 10^{-16}$
MATH	CKA (linear)	0.999	0.635	[0.489, 0.781]	$1.83 \times 10^{-17}$
	unbiased CKA (linear)	1.006	0.629	[0.485, 0.774]	$1.35 \times 10^{-17}$
	CKA (RBF)	0.956	0.674	[0.517, 0.832]	$4.69 \times 10^{-17}$
	unbiased CKA (RBF)	0.970	0.662	[0.508, 0.816]	$3.13 \times 10^{-17}$
TruthfulQA	CKA (linear)	1.012	0.612	[0.470, 0.754]	$2.98 \times 10^{-17}$
	unbiased CKA (linear)	1.023	0.604	[0.466, 0.743]	$1.34 \times 10^{-17}$
	CKA (RBF)	1.005	0.606	[0.460, 0.753]	$4.83 \times 10^{-16}$
	unbiased CKA (RBF)	1.026	0.592	[0.452, 0.732]	$1.37 \times 10^{-16}$

Table 25: Vacation brainstorming task. We fit a mixed-effects regression between mutual information and representational similarity. Here, the summary similarity score is calculated using the average of maximum-aligned scores.  $\beta$  indicates the coefficient of similarity, and the 95% CI represents the confidence interval of  $\beta$ . **The table shows that representational similarity has a significant positive effect on mutual information.**

Dataset	Metric	Intercept	$\beta$	95% CI	$p$
WikiText	CKA (linear)	0.575	0.122	[0.056, 0.188]	0.00031
	unbiased CKA (linear)	0.577	0.121	[0.055, 0.186]	0.00031
	CKA (RBF)	0.586	0.110	[0.045, 0.175]	0.00090
	unbiased CKA (RBF)	0.588	0.108	[0.044, 0.171]	0.00097
GSM8K	CKA (linear)	0.626	0.079	[0.039, 0.119]	$9.26 \times 10^{-5}$
	unbiased CKA (linear)	0.628	0.077	[0.038, 0.116]	0.00011
	CKA (RBF)	0.623	0.082	[0.042, 0.123]	$7.27 \times 10^{-5}$
	unbiased CKA (RBF)	0.626	0.079	[0.040, 0.119]	$8.69 \times 10^{-5}$
MATH	CKA (linear)	0.604	0.105	[0.057, 0.152]	$1.38 \times 10^{-5}$
	unbiased CKA (linear)	0.605	0.103	[0.057, 0.150]	$1.37 \times 10^{-5}$
	CKA (RBF)	0.597	0.111	[0.061, 0.161]	$1.60 \times 10^{-5}$
	unbiased CKA (RBF)	0.599	0.109	[0.060, 0.158]	$1.45 \times 10^{-5}$
TruthfulQA	CKA (linear)	0.613	0.089	[0.043, 0.134]	0.00012
	unbiased CKA (linear)	0.616	0.086	[0.042, 0.130]	0.00014
	CKA (RBF)	0.612	0.089	[0.042, 0.135]	0.00017
	unbiased CKA (RBF)	0.616	0.085	[0.040, 0.129]	0.00018

## E.4 MIXED-EFFECTS

REGRESSION RESULTS FOR MUTUAL INFORMATION CALCULATED WITH LLaMA-3.1-8B

Tables 26~33 show significant positive trends across all datasets, CKA variants, and games.

## E.4.1 WHEN USING THE GLOBAL AVERAGE

Table 26: Story writing task. We fit a mixed-effects regression between mutual information and representational similarity. Here, the summary similarity score is calculated using a global average.  $\beta$  indicates the coefficient of similarity, and the 95% CI represents the confidence interval of  $\beta$ . **The table shows that representational similarity has a significant positive effect on mutual information.**

Dataset	Metric	Intercept	$\beta$	95% CI	$p$
WikiText	CKA (linear)	0.363	0.088	[0.037, 0.139]	0.00075
	unbiased CKA (linear)	0.363	0.087	[0.037, 0.138]	0.00074
	CKA (RBF)	0.372	0.072	[0.026, 0.119]	0.00238
	unbiased CKA (RBF)	0.374	0.071	[0.026, 0.117]	0.00227
GSM8K	CKA (linear)	0.389	0.081	[0.045, 0.118]	$1.39 \times 10^{-5}$
	unbiased CKA (linear)	0.390	0.081	[0.045, 0.117]	$1.20 \times 10^{-5}$
	CKA (RBF)	0.389	0.077	[0.040, 0.113]	$3.60 \times 10^{-5}$
	unbiased CKA (RBF)	0.390	0.077	[0.041, 0.113]	$2.76 \times 10^{-5}$
MATH	CKA (linear)	0.381	0.078	[0.041, 0.116]	$4.25 \times 10^{-5}$
	unbiased CKA (linear)	0.382	0.078	[0.041, 0.115]	$3.90 \times 10^{-5}$
	CKA (RBF)	0.378	0.079	[0.040, 0.118]	$6.24 \times 10^{-5}$
	unbiased CKA (RBF)	0.380	0.078	[0.040, 0.116]	$5.45 \times 10^{-5}$
TruthfulQA	CKA (linear)	0.382	0.079	[0.040, 0.118]	$6.21 \times 10^{-5}$
	unbiased CKA (linear)	0.383	0.080	[0.042, 0.119]	$4.33 \times 10^{-5}$
	CKA (RBF)	0.383	0.071	[0.033, 0.110]	0.00030
	unbiased CKA (RBF)	0.385	0.072	[0.035, 0.110]	0.00018

Table 27: Fictional biography generation task. We fit a mixed-effects regression between mutual information and representational similarity. Here, the summary similarity score is calculated using a global average.  $\beta$  indicates the coefficient of similarity, and the 95% CI represents the confidence interval of  $\beta$ . **The table shows that representational similarity has a significant positive effect on mutual information.**

Dataset	Metric	Intercept	$\beta$	95% CI	$p$
WikiText	CKA (linear)	0.135	0.130	[0.098, 0.162]	$1.88 \times 10^{-15}$
	unbiased CKA (linear)	0.136	0.129	[0.097, 0.161]	$1.70 \times 10^{-15}$
	CKA (RBF)	0.136	0.130	[0.0981, 0.163]	$2.96 \times 10^{-15}$
	unbiased CKA (RBF)	0.138	0.128	[0.097, 0.160]	$2.12 \times 10^{-15}$
GSM8K	CKA (linear)	0.175	0.119	[0.094, 0.144]	$5.79 \times 10^{-21}$
	unbiased CKA (linear)	0.177	0.118	[0.093, 0.142]	$5.51 \times 10^{-21}$
	CKA (RBF)	0.173	0.118	[0.093, 0.143]	$1.74 \times 10^{-20}$
	unbiased CKA (RBF)	0.176	0.116	[0.091, 0.140]	$1.37 \times 10^{-20}$
MATH	CKA (linear)	0.162	0.118	[0.093, 0.144]	$7.24 \times 10^{-20}$
	unbiased CKA (linear)	0.164	0.117	[0.092, 0.142]	$6.08 \times 10^{-20}$
	CKA (RBF)	0.155	0.123	[0.097, 0.150]	$9.91 \times 10^{-20}$
	unbiased CKA (RBF)	0.158	0.121	[0.095, 0.147]	$7.60 \times 10^{-20}$
TruthfulQA	CKA (linear)	0.169	0.106	[0.079, 0.132]	$5.53 \times 10^{-15}$
	unbiased CKA (linear)	0.170	0.106	[0.080, 0.132]	$2.24 \times 10^{-15}$
	CKA (RBF)	0.166	0.105	[0.078, 0.132]	$3.05 \times 10^{-14}$
	unbiased CKA (RBF)	0.169	0.104	[0.078, 0.130]	$7.51 \times 10^{-15}$

3078  
 3079  
 3080  
 3081  
 3082  
 3083  
 3084  
 3085  
 3086  
 3087  
 3088  
 3089  
 3090  
 3091  
 3092  
 3093  
 3094  
 3095  
 3096  
 3097  
 3098  
 3099  
 3100  
 3101  
 3102  
 3103  
 3104  
 3105  
 3106  
 3107  
 3108  
 3109  
 3110  
 3111  
 3112  
 3113  
 3114  
 3115  
 3116  
 3117  
 3118  
 3119  
 3120  
 3121  
 3122  
 3123  
 3124  
 3125  
 3126  
 3127  
 3128  
 3129  
 3130  
 3131

Table 28: Haiku composition task. We fit a mixed-effects regression between mutual information and representational similarity. Here, the summary similarity score is calculated using a global average.  $\beta$  indicates the coefficient of similarity, and the 95% CI represents the confidence interval of  $\beta$ . **The table shows that representational similarity has a significant positive effect on mutual information.**

Dataset	Metric	Intercept	$\beta$	95% CI	$p$
WikiText	CKA (linear)	0.628	1.335	[1.033, 1.638]	$5.17 \times 10^{-18}$
	unbiased CKA (linear)	0.637	1.329	[1.029, 1.629]	$3.54 \times 10^{-18}$
	CKA (RBF)	0.700	1.230	[0.928, 1.533]	$1.71 \times 10^{-15}$
	unbiased CKA (RBF)	0.720	1.215	[0.918, 1.511]	$1.02 \times 10^{-15}$
GSM8K	CKA (linear)	1.209	0.693	[0.466, 0.920]	$2.27 \times 10^{-9}$
	unbiased CKA (linear)	1.218	0.684	[0.460, 0.909]	$2.23 \times 10^{-9}$
	CKA (RBF)	1.200	0.675	[0.447, 0.903]	$6.47 \times 10^{-9}$
	unbiased CKA (RBF)	1.216	0.662	[0.439, 0.884]	$5.65 \times 10^{-9}$
MATH	CKA (linear)	1.058	0.863	[0.628, 1.098]	$6.30 \times 10^{-13}$
	unbiased CKA (linear)	1.066	0.861	[0.629, 1.093]	$3.39 \times 10^{-13}$
	CKA (RBF)	1.011	0.897	[0.651, 1.143]	$8.97 \times 10^{-13}$
	unbiased CKA (RBF)	1.026	0.891	[0.651, 1.131]	$3.44 \times 10^{-13}$
TruthfulQA	CKA (linear)	1.098	0.791	[0.545, 1.036]	$2.65 \times 10^{-10}$
	unbiased CKA (linear)	1.103	0.804	[0.562, 1.046]	$7.32 \times 10^{-11}$
	CKA (RBF)	1.090	0.762	[0.512, 1.012]	$2.22 \times 10^{-9}$
	unbiased CKA (RBF)	1.103	0.777	[0.535, 1.020]	$3.34 \times 10^{-10}$

Table 29: Vacation brainstorming task. We fit a mixed-effects regression between mutual information and representational similarity. Here, the summary similarity score is calculated using a global average.  $\beta$  indicates the coefficient of similarity, and the 95% CI represents the confidence interval of  $\beta$ . **The table shows that representational similarity has a significant positive effect on mutual information.**

Dataset	Metric	Intercept	$\beta$	95% CI	$p$
WikiText	CKA (linear)	0.696	0.158	[0.094, 0.222]	$1.40 \times 10^{-6}$
	unbiased CKA (linear)	0.698	0.155	[0.094, 0.215]	$6.16 \times 10^{-7}$
	CKA (RBF)	0.698	0.156	[0.087, 0.224]	$8.20 \times 10^{-6}$
	unbiased CKA (RBF)	0.701	0.153	[0.086, 0.220]	$7.34 \times 10^{-6}$
GSM8K	CKA (linear)	0.733	0.177	[0.112, 0.243]	$1.20 \times 10^{-7}$
	unbiased CKA (linear)	0.735	0.177	[0.112, 0.242]	$9.84 \times 10^{-8}$
	CKA (RBF)	0.730	0.175	[0.110, 0.241]	$1.83 \times 10^{-7}$
	unbiased CKA (RBF)	0.734	0.175	[0.110, 0.239]	$1.25 \times 10^{-7}$
MATH	CKA (linear)	0.713	0.179	[0.117, 0.241]	$1.39 \times 10^{-8}$
	unbiased CKA (linear)	0.715	0.179	[0.118, 0.240]	$1.01 \times 10^{-8}$
	CKA (RBF)	0.703	0.186	[0.121, 0.251]	$1.84 \times 10^{-8}$
	unbiased CKA (RBF)	0.707	0.184	[0.121, 0.247]	$1.13 \times 10^{-8}$
TruthfulQA	CKA (linear)	0.731	0.142	[0.078, 0.206]	$1.33 \times 10^{-5}$
	unbiased CKA (linear)	0.733	0.142	[0.079, 0.206]	$1.08 \times 10^{-5}$
	CKA (RBF)	0.729	0.138	[0.073, 0.203]	$3.04 \times 10^{-5}$
	unbiased CKA (RBF)	0.733	0.137	[0.074, 0.200]	$2.09 \times 10^{-5}$

## E.4.2 WHEN USING THE AVERAGE OF MAXIMUM-ALIGNED SCORES

Table 30: Story writing task. We fit a mixed-effects regression between mutual information and representational similarity. Here, the summary similarity score is calculated using the average of maximum-aligned scores.  $\beta$  indicates the coefficient of similarity, and the 95% CI represents the confidence interval of  $\beta$ . **The table shows that representational similarity has a significant positive effect on mutual information.**

Dataset	Metric	Intercept	$\beta$	95% CI	$p$
WikiText	CKA (linear)	0.317	0.131	[0.090, 0.171]	$2.92 \times 10^{-10}$
	unbiased CKA (linear)	0.318	0.129	[0.089, 0.169]	$3.11 \times 10^{-10}$
	CKA (RBF)	0.327	0.120	[0.081, 0.159]	$2.00 \times 10^{-9}$
	unbiased CKA (RBF)	0.329	0.118	[0.079, 0.156]	$1.96 \times 10^{-9}$
GSM8K	CKA (linear)	0.380	0.069	[0.046, 0.092]	$2.98 \times 10^{-9}$
	unbiased CKA (linear)	0.381	0.068	[0.046, 0.091]	$2.99 \times 10^{-9}$
	CKA (RBF)	0.378	0.070	[0.046, 0.093]	$6.20 \times 10^{-9}$
	unbiased CKA (RBF)	0.380	0.068	[0.045, 0.091]	$5.88 \times 10^{-9}$
MATH	CKA (linear)	0.362	0.088	[0.060, 0.115]	$4.06 \times 10^{-10}$
	unbiased CKA (linear)	0.364	0.086	[0.059, 0.113]	$4.40 \times 10^{-10}$
	CKA (RBF)	0.357	0.092	[0.063, 0.121]	$8.61 \times 10^{-10}$
	unbiased CKA (RBF)	0.359	0.090	[0.061, 0.118]	$9.81 \times 10^{-10}$
TruthfulQA	CKA (linear)	0.363	0.086	[0.059, 0.113]	$2.37 \times 10^{-10}$
	unbiased CKA (linear)	0.365	0.084	[0.058, 0.110]	$2.25 \times 10^{-10}$
	CKA (RBF)	0.363	0.085	[0.057, 0.112]	$1.46 \times 10^{-9}$
	unbiased CKA (RBF)	0.366	0.082	[0.055, 0.108]	$1.16 \times 10^{-9}$

Table 31: Fictional biography generation task. We fit a mixed-effects regression between mutual information and representational similarity. Here, the summary similarity score is calculated using the average of maximum-aligned scores.  $\beta$  indicates the coefficient of similarity, and the 95% CI represents the confidence interval of  $\beta$ . **The table shows that representational similarity has a significant positive effect on mutual information.**

Dataset	Metric	Intercept	$\beta$	95% CI	$p$
WikiText	CKA (linear)	0.143	0.094	[0.068, 0.120]	$1.60 \times 10^{-12}$
	unbiased CKA (linear)	0.144	0.093	[0.067, 0.119]	$1.83 \times 10^{-12}$
	CKA (RBF)	0.144	0.095	[0.069, 0.120]	$2.72 \times 10^{-13}$
	unbiased CKA (RBF)	0.146	0.092	[0.068, 0.117]	$3.19 \times 10^{-13}$
GSM8K	CKA (linear)	0.178	0.069	[0.054, 0.084]	$5.27 \times 10^{-20}$
	unbiased CKA (linear)	0.179	0.068	[0.054, 0.083]	$6.93 \times 10^{-20}$
	CKA (RBF)	0.175	0.072	[0.057, 0.088]	$1.65 \times 10^{-20}$
	unbiased CKA (RBF)	0.177	0.070	[0.055, 0.085]	$2.42 \times 10^{-20}$
MATH	CKA (linear)	0.167	0.077	[0.059, 0.095]	$3.49 \times 10^{-17}$
	unbiased CKA (linear)	0.168	0.075	[0.058, 0.093]	$5.06 \times 10^{-17}$
	CKA (RBF)	0.161	0.083	[0.064, 0.103]	$1.85 \times 10^{-17}$
	unbiased CKA (RBF)	0.163	0.081	[0.062, 0.100]	$3.36 \times 10^{-17}$
TruthfulQA	CKA (linear)	0.170	0.071	[0.054, 0.089]	$6.60 \times 10^{-16}$
	unbiased CKA (linear)	0.172	0.070	[0.053, 0.086]	$7.44 \times 10^{-16}$
	CKA (RBF)	0.168	0.072	[0.054, 0.090]	$1.93 \times 10^{-15}$
	unbiased CKA (RBF)	0.172	0.069	[0.052, 0.086]	$1.85 \times 10^{-15}$

3186  
3187  
3188  
3189  
3190  
3191  
3192  
3193  
3194  
3195  
3196  
3197  
3198  
3199  
3200  
3201  
3202  
3203  
3204  
3205  
3206  
3207  
3208  
3209  
3210  
3211  
3212  
3213  
3214  
3215  
3216  
3217  
3218  
3219  
3220  
3221  
3222  
3223  
3224  
3225  
3226  
3227  
3228  
3229  
3230  
3231  
3232  
3233  
3234  
3235  
3236  
3237  
3238  
3239

Table 32: Haiku composition task. We fit a mixed-effects regression between mutual information and representational similarity. Here, the summary similarity score is calculated using the average of maximum-aligned scores.  $\beta$  indicates the coefficient of similarity, and the 95% CI represents the confidence interval of  $\beta$ . **The table shows that representational similarity has a significant positive effect on mutual information.**

Dataset	Metric	Intercept	$\beta$	95% CI	$p$
WikiText	CKA (linear)	0.664	1.023	[0.788, 1.259]	$1.73 \times 10^{-17}$
	unbiased CKA (linear)	0.672	1.015	[0.782, 1.249]	$1.55 \times 10^{-17}$
	CKA (RBF)	0.726	0.957	[0.728, 1.187]	$2.83 \times 10^{-16}$
	unbiased CKA (RBF)	0.743	0.942	[0.717, 1.167]	$2.20 \times 10^{-16}$
GSM8K	CKA (linear)	1.161	0.530	[0.395, 0.664]	$9.96 \times 10^{-15}$
	unbiased CKA (linear)	1.169	0.523	[0.391, 0.655]	$8.96 \times 10^{-15}$
	CKA (RBF)	1.149	0.535	[0.397, 0.673]	$3.24 \times 10^{-14}$
	unbiased CKA (RBF)	1.164	0.522	[0.387, 0.656]	$3.02 \times 10^{-14}$
MATH	CKA (linear)	1.031	0.664	[0.503, 0.826]	$7.19 \times 10^{-16}$
	unbiased CKA (linear)	1.039	0.658	[0.499, 0.817]	$5.55 \times 10^{-16}$
	CKA (RBF)	0.985	0.706	[0.532, 0.879]	$1.57 \times 10^{-15}$
	unbiased CKA (RBF)	1.000	0.693	[0.523, 0.862]	$1.14 \times 10^{-15}$
TruthfulQA	CKA (linear)	1.049	0.632	[0.475, 0.788]	$2.49 \times 10^{-15}$
	unbiased CKA (linear)	1.061	0.624	[0.471, 0.777]	$1.22 \times 10^{-15}$
	CKA (RBF)	1.042	0.628	[0.466, 0.789]	$2.55 \times 10^{-14}$
	unbiased CKA (RBF)	1.063	0.613	[0.458, 0.768]	$8.33 \times 10^{-15}$

Table 33: Vacation brainstorming task. We fit a mixed-effects regression between mutual information and representational similarity. Here, the summary similarity score is calculated using the average of maximum-aligned scores.  $\beta$  indicates the coefficient of similarity, and the 95% CI represents the confidence interval of  $\beta$ . **The table shows that representational similarity has a significant positive effect on mutual information.**

Dataset	Metric	Intercept	$\beta$	95% CI	$p$
WikiText	CKA (linear)	0.654	0.182	[0.114, 0.249]	$1.36 \times 10^{-7}$
	unbiased CKA (linear)	0.656	0.180	[0.113, 0.247]	$1.40 \times 10^{-7}$
	CKA (RBF)	0.666	0.169	[0.102, 0.236]	$7.80 \times 10^{-7}$
	unbiased CKA (RBF)	0.670	0.165	[0.099, 0.231]	$8.44 \times 10^{-7}$
GSM8K	CKA (linear)	0.731	0.117	[0.076, 0.159]	$3.98 \times 10^{-8}$
	unbiased CKA (linear)	0.733	0.115	[0.074, 0.157]	$4.77 \times 10^{-8}$
	CKA (RBF)	0.726	0.121	[0.078, 0.165]	$3.29 \times 10^{-8}$
	unbiased CKA (RBF)	0.730	0.118	[0.075, 0.160]	$4.21 \times 10^{-8}$
MATH	CKA (linear)	0.701	0.148	[0.098, 0.197]	$4.32 \times 10^{-9}$
	unbiased CKA (linear)	0.703	0.146	[0.097, 0.195]	$4.16 \times 10^{-9}$
	CKA (RBF)	0.691	0.158	[0.105, 0.211]	$4.75 \times 10^{-9}$
	unbiased CKA (RBF)	0.694	0.155	[0.103, 0.207]	$4.03 \times 10^{-9}$
TruthfulQA	CKA (linear)	0.712	0.131	[0.083, 0.178]	$5.71 \times 10^{-8}$
	unbiased CKA (linear)	0.715	0.127	[0.081, 0.173]	$6.40 \times 10^{-8}$
	CKA (RBF)	0.709	0.132	[0.083, 0.180]	$8.78 \times 10^{-8}$
	unbiased CKA (RBF)	0.715	0.126	[0.080, 0.173]	$8.95 \times 10^{-8}$

## E.5 MIXED-EFFECTS

## REGRESSION RESULTS FOR MUTUAL INFORMATION, EXCLUDING THE LLAMA FAMILY

Tables 34~41 show significant positive trends across all datasets, CKA variants, and games.

## E.5.1 WHEN USING THE GLOBAL AVERAGE

Table 34: Story writing task. We fit a mixed-effects regression between mutual information and representational similarity. Here, the summary similarity score is calculated using a global average.  $\beta$  indicates the coefficient of similarity, and the 95% CI represents the confidence interval of  $\beta$ . **The table shows that representational similarity has a significant positive effect on mutual information.**

Dataset	Metric	Intercept	$\beta$	95% CI	$p$
WikiText	CKA (linear)	0.307	0.058	[0.020, 0.095]	0.00246
	unbiased CKA (linear)	0.308	0.058	[0.021, 0.095]	0.00232
	CKA (RBF)	0.311	0.052	[0.015, 0.090]	0.00633
	unbiased CKA (RBF)	0.312	0.052	[0.015, 0.089]	0.00584
GSM8K	CKA (linear)	0.323	0.058	[0.024, 0.092]	0.00083
	unbiased CKA (linear)	0.324	0.058	[0.024, 0.092]	0.00076
	CKA (RBF)	0.323	0.056	[0.022, 0.090]	0.00130
	unbiased CKA (RBF)	0.324	0.056	[0.022, 0.089]	0.00113
MATH	CKA (linear)	0.320	0.051	[0.018, 0.083]	0.00212
	unbiased CKA (linear)	0.320	0.051	[0.019, 0.083]	0.00192
	CKA (RBF)	0.317	0.053	[0.019, 0.086]	0.00193
	unbiased CKA (RBF)	0.318	0.053	[0.020, 0.085]	0.00165
TruthfulQA	CKA (linear)	0.318	0.056	[0.023, 0.090]	0.00097
	unbiased CKA (linear)	0.319	0.057	[0.024, 0.090]	0.00067
	CKA (RBF)	0.319	0.053	[0.020, 0.087]	0.00186
	unbiased CKA (RBF)	0.319	0.054	[0.022, 0.087]	0.00116

Table 35: Fictional biography generation task. We fit a mixed-effects regression between mutual information and representational similarity. Here, the summary similarity score is calculated using a global average.  $\beta$  indicates the coefficient of similarity, and the 95% CI represents the confidence interval of  $\beta$ . **The table shows that representational similarity has a significant positive effect on mutual information.**

Dataset	Metric	Intercept	$\beta$	95% CI	$p$
WikiText	CKA (linear)	0.078	0.138	[0.105, 0.171]	$5.0 \times 10^{-16}$
	unbiased CKA (linear)	0.080	0.136	[0.103, 0.169]	$4.8 \times 10^{-16}$
	CKA (RBF)	0.081	0.135	[0.102, 0.169]	$2.2 \times 10^{-15}$
	unbiased CKA (RBF)	0.084	0.133	[0.100, 0.166]	$1.9 \times 10^{-15}$
GSM8K	CKA (linear)	0.121	0.121	[0.095, 0.148]	$6.0 \times 10^{-19}$
	unbiased CKA (linear)	0.123	0.120	[0.094, 0.147]	$5.0 \times 10^{-19}$
	CKA (RBF)	0.119	0.120	[0.093, 0.147]	$1.7 \times 10^{-18}$
	unbiased CKA (RBF)	0.122	0.118	[0.092, 0.145]	$1.1 \times 10^{-18}$
MATH	CKA (linear)	0.108	0.120	[0.094, 0.146]	$3.2 \times 10^{-19}$
	unbiased CKA (linear)	0.109	0.119	[0.093, 0.145]	$2.5 \times 10^{-19}$
	CKA (RBF)	0.102	0.124	[0.097, 0.152]	$1.0 \times 10^{-18}$
	unbiased CKA (RBF)	0.104	0.122	[0.095, 0.149]	$6.8 \times 10^{-19}$
TruthfulQA	CKA (linear)	0.116	0.107	[0.079, 0.135]	$4.3 \times 10^{-14}$
	unbiased CKA (linear)	0.117	0.107	[0.079, 0.134]	$2.3 \times 10^{-14}$
	CKA (RBF)	0.115	0.103	[0.075, 0.132]	$6.4 \times 10^{-13}$
	unbiased CKA (RBF)	0.118	0.103	[0.075, 0.130]	$2.2 \times 10^{-13}$

3294  
3295  
3296  
3297

Table 36: Haiku composition task. We fit a mixed-effects regression between mutual information and representational similarity. Here, the summary similarity score is calculated using a global average.  $\beta$  indicates the coefficient of similarity, and the 95% CI represents the confidence interval of  $\beta$ . **The table shows that representational similarity has a significant positive effect on mutual information.**

3301  
3302  
3303  
3304  
3305  
3306  
3307  
3308  
3309  
3310  
3311  
3312  
3313  
3314  
3315  
3316  
3317

Dataset	Metric	Intercept	$\beta$	95% CI	$p$
WikiText	CKA (linear)	0.497	1.536	[1.242, 1.831]	$1.7 \times 10^{-24}$
	unbiased CKA (linear)	0.510	1.522	[1.231, 1.814]	$1.4 \times 10^{-24}$
	CKA (RBF)	0.590	1.404	[1.108, 1.700]	$1.5 \times 10^{-20}$
	unbiased CKA (RBF)	0.615	1.381	[1.092, 1.670]	$7.1 \times 10^{-21}$
GSM8K	CKA (linear)	1.106	0.958	[0.729, 1.188]	$2.6 \times 10^{-16}$
	unbiased CKA (linear)	1.117	0.951	[0.724, 1.178]	$2.2 \times 10^{-16}$
	CKA (RBF)	1.092	0.941	[0.712, 1.171]	$9.8 \times 10^{-16}$
	unbiased CKA (RBF)	1.111	0.929	[0.704, 1.154]	$6.4 \times 10^{-16}$
MATH	CKA (linear)	0.978	1.001	[0.772, 1.230]	$1.0 \times 10^{-17}$
	unbiased CKA (linear)	0.987	0.996	[0.770, 1.222]	$5.8 \times 10^{-18}$
	CKA (RBF)	0.913	1.061	[0.820, 1.302]	$5.5 \times 10^{-18}$
	unbiased CKA (RBF)	0.932	1.048	[0.813, 1.283]	$2.2 \times 10^{-18}$
TruthfulQA	CKA (linear)	0.991	1.014	[0.772, 1.256]	$2.1 \times 10^{-16}$
	unbiased CKA (linear)	1.003	1.018	[0.780, 1.257]	$5.8 \times 10^{-17}$
	CKA (RBF)	0.987	0.972	[0.727, 1.217]	$7.6 \times 10^{-15}$
	unbiased CKA (RBF)	1.008	0.975	[0.737, 1.213]	$9.7 \times 10^{-16}$

3318  
3319  
3320  
3321  
3322  
3323  
3324

Table 37: Vacation brainstorming task. We fit a mixed-effects regression between mutual information and representational similarity. Here, the summary similarity score is calculated using a global average.  $\beta$  indicates the coefficient of similarity, and the 95% CI represents the confidence interval of  $\beta$ . **The table shows that representational similarity has a significant positive effect on mutual information.**

3325  
3326  
3327  
3328  
3329  
3330  
3331  
3332  
3333  
3334  
3335  
3336  
3337  
3338  
3339  
3340  
3341  
3342  
3343  
3344  
3345

Dataset	Metric	Intercept	$\beta$	95% CI	$p$
WikiText	CKA (linear)	0.592	0.185	[0.119, 0.250]	$3.2 \times 10^{-8}$
	unbiased CKA (linear)	0.593	0.183	[0.118, 0.248]	$3.2 \times 10^{-8}$
	CKA (RBF)	0.597	0.179	[0.107, 0.251]	$1.1 \times 10^{-6}$
	unbiased CKA (RBF)	0.600	0.176	[0.105, 0.246]	$1.0 \times 10^{-6}$
GSM8K	CKA (linear)	0.644	0.180	[0.108, 0.251]	$7.8 \times 10^{-7}$
	unbiased CKA (linear)	0.646	0.179	[0.108, 0.250]	$7.3 \times 10^{-7}$
	CKA (RBF)	0.642	0.176	[0.105, 0.247]	$1.2 \times 10^{-6}$
	unbiased CKA (RBF)	0.645	0.175	[0.105, 0.246]	$9.7 \times 10^{-7}$
MATH	CKA (linear)	0.623	0.182	[0.117, 0.247]	$4.3 \times 10^{-8}$
	unbiased CKA (linear)	0.624	0.181	[0.117, 0.245]	$3.3 \times 10^{-8}$
	CKA (RBF)	0.613	0.188	[0.120, 0.255]	$4.5 \times 10^{-8}$
	unbiased CKA (RBF)	0.616	0.186	[0.120, 0.252]	$2.8 \times 10^{-8}$
TruthfulQA	CKA (linear)	0.638	0.154	[0.083, 0.224]	$1.9 \times 10^{-5}$
	unbiased CKA (linear)	0.640	0.153	[0.084, 0.223]	$1.7 \times 10^{-5}$
	CKA (RBF)	0.638	0.146	[0.075, 0.216]	$5.6 \times 10^{-5}$
	unbiased CKA (RBF)	0.642	0.145	[0.076, 0.214]	$4.1 \times 10^{-5}$

3346  
3347

## E.5.2 WHEN USING THE AVERAGE OF MAXIMUM-ALIGNED SCORES

Table 38: Story writing task. We fit a mixed-effects regression between mutual information and representational similarity. Here, the summary similarity score is calculated using the average of maximum-aligned scores.  $\beta$  indicates the coefficient of similarity, and the 95% CI represents the confidence interval of  $\beta$ . **The table shows that representational similarity has a significant positive effect on mutual information.**

Dataset	Metric	Intercept	$\beta$	95% CI	$p$
WikiText	CKA (linear)	0.275	0.090	[0.054, 0.125]	$5.87 \times 10^{-7}$
	unbiased CKA (linear)	0.275	0.089	[0.054, 0.124]	$5.71 \times 10^{-7}$
	CKA (RBF)	0.281	0.083	[0.049, 0.118]	$1.94 \times 10^{-6}$
	unbiased CKA (RBF)	0.282	0.082	[0.048, 0.116]	$1.81 \times 10^{-6}$
GSM8K	CKA (linear)	0.318	0.047	[0.026, 0.068]	$1.40 \times 10^{-5}$
	unbiased CKA (linear)	0.318	0.047	[0.026, 0.068]	$1.36 \times 10^{-5}$
	CKA (RBF)	0.317	0.048	[0.026, 0.069]	$1.99 \times 10^{-5}$
	unbiased CKA (RBF)	0.318	0.047	[0.025, 0.068]	$1.89 \times 10^{-5}$
MATH	CKA (linear)	0.306	0.059	[0.034, 0.084]	$3.48 \times 10^{-6}$
	unbiased CKA (linear)	0.307	0.059	[0.034, 0.083]	$3.46 \times 10^{-6}$
	CKA (RBF)	0.302	0.063	[0.036, 0.090]	$4.00 \times 10^{-6}$
	unbiased CKA (RBF)	0.303	0.062	[0.035, 0.088]	$3.92 \times 10^{-6}$
TruthfulQA	CKA (linear)	0.307	0.058	[0.034, 0.082]	$1.75 \times 10^{-6}$
	unbiased CKA (linear)	0.308	0.058	[0.034, 0.081]	$1.43 \times 10^{-6}$
	CKA (RBF)	0.306	0.058	[0.033, 0.082]	$3.66 \times 10^{-6}$
	unbiased CKA (RBF)	0.308	0.056	[0.033, 0.080]	$2.63 \times 10^{-6}$

Table 39: Fictional biography generation task. We fit a mixed-effects regression between mutual information and representational similarity. Here, the summary similarity score is calculated using the average of maximum-aligned scores.  $\beta$  indicates the coefficient of similarity, and the 95% CI represents the confidence interval of  $\beta$ . **The table shows that representational similarity has a significant positive effect on mutual information.**

Dataset	Metric	Intercept	$\beta$	95% CI	$p$
WikiText	CKA (linear)	0.086	0.099	[0.072, 0.126]	$3.39 \times 10^{-13}$
	unbiased CKA (linear)	0.087	0.098	[0.072, 0.124]	$3.67 \times 10^{-13}$
	CKA (RBF)	0.089	0.097	[0.071, 0.123]	$1.93 \times 10^{-13}$
	unbiased CKA (RBF)	0.091	0.095	[0.070, 0.121]	$2.12 \times 10^{-13}$
GSM8K	CKA (linear)	0.124	0.070	[0.055, 0.086]	$4.29 \times 10^{-19}$
	unbiased CKA (linear)	0.125	0.069	[0.054, 0.085]	$5.20 \times 10^{-19}$
	CKA (RBF)	0.121	0.073	[0.057, 0.089]	$1.94 \times 10^{-19}$
	unbiased CKA (RBF)	0.123	0.071	[0.056, 0.087]	$2.50 \times 10^{-19}$
MATH	CKA (linear)	0.112	0.080	[0.062, 0.098]	$1.85 \times 10^{-17}$
	unbiased CKA (linear)	0.113	0.079	[0.061, 0.097]	$2.34 \times 10^{-17}$
	CKA (RBF)	0.106	0.086	[0.066, 0.106]	$2.03 \times 10^{-17}$
	unbiased CKA (RBF)	0.108	0.084	[0.064, 0.103]	$3.01 \times 10^{-17}$
TruthfulQA	CKA (linear)	0.116	0.073	[0.055, 0.091]	$9.57 \times 10^{-16}$
	unbiased CKA (linear)	0.118	0.071	[0.054, 0.089]	$9.84 \times 10^{-16}$
	CKA (RBF)	0.115	0.073	[0.055, 0.091]	$4.75 \times 10^{-15}$
	unbiased CKA (RBF)	0.118	0.070	[0.053, 0.088]	$4.28 \times 10^{-15}$

3402  
3403  
3404  
3405  
3406  
3407  
3408  
3409  
3410  
3411  
3412  
3413  
3414  
3415  
3416  
3417  
3418  
3419  
3420  
3421  
3422  
3423  
3424  
3425  
3426  
3427  
3428  
3429  
3430  
3431  
3432  
3433  
3434  
3435  
3436  
3437  
3438  
3439  
3440  
3441  
3442  
3443  
3444  
3445  
3446  
3447  
3448  
3449  
3450  
3451  
3452  
3453  
3454  
3455

Table 40: Haiku composition task. We fit a mixed-effects regression between mutual information and representational similarity. Here, the summary similarity score is calculated using the average of maximum-aligned scores.  $\beta$  indicates the coefficient of similarity, and the 95% CI represents the confidence interval of  $\beta$ . **The table shows that representational similarity has a significant positive effect on mutual information.**

Dataset	Metric	Intercept	$\beta$	95% CI	$p$
WikiText	CKA (linear)	0.628	1.048	[0.823, 1.274]	$8.91 \times 10^{-20}$
	unbiased CKA (linear)	0.637	1.040	[0.816, 1.264]	$8.54 \times 10^{-20}$
	CKA (RBF)	0.702	0.970	[0.750, 1.189]	$4.29 \times 10^{-18}$
	unbiased CKA (RBF)	0.719	0.953	[0.738, 1.168]	$3.67 \times 10^{-18}$
GSM8K	CKA (linear)	1.104	0.602	[0.472, 0.733]	$1.22 \times 10^{-19}$
	unbiased CKA (linear)	1.112	0.595	[0.466, 0.723]	$1.24 \times 10^{-19}$
	CKA (RBF)	1.088	0.613	[0.479, 0.747]	$3.39 \times 10^{-19}$
	unbiased CKA (RBF)	1.103	0.599	[0.468, 0.730]	$3.27 \times 10^{-19}$
MATH	CKA (linear)	0.984	0.710	[0.554, 0.865]	$3.73 \times 10^{-19}$
	unbiased CKA (linear)	0.992	0.701	[0.547, 0.855]	$3.93 \times 10^{-19}$
	CKA (RBF)	0.934	0.756	[0.589, 0.923]	$8.06 \times 10^{-19}$
	unbiased CKA (RBF)	0.951	0.739	[0.575, 0.903]	$9.36 \times 10^{-19}$
TruthfulQA	CKA (linear)	0.990	0.701	[0.551, 0.852]	$6.29 \times 10^{-20}$
	unbiased CKA (linear)	1.006	0.686	[0.539, 0.833]	$6.31 \times 10^{-20}$
	CKA (RBF)	0.986	0.693	[0.538, 0.847]	$1.48 \times 10^{-18}$
	unbiased CKA (RBF)	1.013	0.669	[0.521, 0.817]	$9.74 \times 10^{-19}$

Table 41: Vacation brainstorming task. We fit a mixed-effects regression between mutual information and representational similarity. Here, the summary similarity score is calculated using the average of maximum-aligned scores.  $\beta$  indicates the coefficient of similarity, and the 95% CI represents the confidence interval of  $\beta$ . **The table shows that representational similarity has a significant positive effect on mutual information.**

Dataset	Metric	Intercept	$\beta$	95% CI	$p$
WikiText	CKA (linear)	0.582	0.161	[0.086, 0.235]	$2.21 \times 10^{-5}$
	unbiased CKA (linear)	0.583	0.159	[0.086, 0.233]	$2.22 \times 10^{-5}$
	CKA (RBF)	0.596	0.144	[0.072, 0.216]	$9.66 \times 10^{-5}$
	unbiased CKA (RBF)	0.599	0.141	[0.070, 0.212]	$9.68 \times 10^{-5}$
GSM8K	CKA (linear)	0.651	0.099	[0.055, 0.142]	$9.04 \times 10^{-6}$
	unbiased CKA (linear)	0.653	0.097	[0.054, 0.140]	$1.01 \times 10^{-5}$
	CKA (RBF)	0.648	0.102	[0.057, 0.147]	$8.23 \times 10^{-6}$
	unbiased CKA (RBF)	0.650	0.099	[0.055, 0.143]	$9.33 \times 10^{-6}$
MATH	CKA (linear)	0.624	0.129	[0.077, 0.181]	$1.13 \times 10^{-6}$
	unbiased CKA (linear)	0.625	0.128	[0.077, 0.180]	$9.97 \times 10^{-7}$
	CKA (RBF)	0.615	0.137	[0.081, 0.193]	$1.34 \times 10^{-6}$
	unbiased CKA (RBF)	0.617	0.136	[0.081, 0.190]	$9.86 \times 10^{-7}$
TruthfulQA	CKA (linear)	0.633	0.114	[0.064, 0.164]	$7.39 \times 10^{-6}$
	unbiased CKA (linear)	0.635	0.112	[0.063, 0.161]	$7.23 \times 10^{-6}$
	CKA (RBF)	0.632	0.113	[0.062, 0.163]	$1.43 \times 10^{-5}$
	unbiased CKA (RBF)	0.637	0.109	[0.060, 0.158]	$1.25 \times 10^{-5}$

## E.6 MIXED-EFFECTS

## REGRESSION RESULTS OBTAINED BY CONTROLLING FOR PERFORMANCE DISPARITY

Table 42: Results of the mixed-effects regression controlling for performance disparity between the two models. Performance disparity is measured as the difference in their MMLU scores. The summary similarity score is computed as the global average representational similarity obtained using linear CKA on WikiText. Here,  $\beta$  denotes the regression coefficient, and the 95% CI refers to the 95% confidence interval of  $\beta$ . **The table shows that representational similarity has a significant effect on cooperation and novelty even when controlling for performance disparity. This implies that the trend is not merely a byproduct of performance disparities.**

Game	Predictor	$\beta$	95% CI	$p$
Word Guessing	Representational Similarity	6.957	[5.160, 8.754]	< 0.001
	Performance Disparity	-0.829	[-1.710, 0.053]	0.066
Public Good	Representational Similarity	48.846	[27.680, 70.012]	< 0.001
	Performance Disparity	-14.418	[-30.133, 1.297]	0.072
Divide-a-Dollar	Representational Similarity	0.452	[0.221, 0.684]	< 0.001
	Performance Disparity	0.032	[-0.110, 0.174]	0.656
KBC	Representational Similarity	13.406	[-0.762, 27.574]	0.064
	Performance Disparity	-22.757	[-33.156, -12.357]	< 0.001
Story (Uniqueness)	Representational Similarity	-0.268	[-1.455, 0.919]	0.658
	Performance Disparity	0.320	[-0.635, 1.275]	0.511
Biography (Uniqueness)	Representational Similarity	-0.696	[-1.814, 0.423]	0.223
	Performance Disparity	1.580	[0.868, 2.293]	< 0.001
Haiku (Uniqueness)	Representational Similarity	-3.488	[-4.810, -2.167]	< 0.001
	Performance Disparity	-0.012	[-1.052, 1.028]	0.982
Vacation (Uniqueness)	Representational Similarity	-1.046	[-1.625, -0.468]	< 0.001
	Performance Disparity	-0.414	[-0.955, 0.128]	0.134
Story (MI)	Representational Similarity	0.054	[0.019, 0.088]	0.002
	Performance Disparity	0.001	[-0.025, 0.028]	0.935
Biography (MI)	Representational Similarity	0.124	[0.094, 0.154]	< 0.001
	Performance Disparity	0.009	[-0.010, 0.027]	0.363
Haiku (MI)	Representational Similarity	1.353	[1.069, 1.637]	< 0.001
	Performance Disparity	0.107	[-0.059, 0.273]	0.208
Vacation (MI)	Representational Similarity	0.130	[0.068, 0.192]	< 0.001
	Performance Disparity	-0.014	[-0.066, 0.038]	0.600

## E.7 MIXED-EFFECTS REGRESSION RESULTS FOR EACH LAYER GROUP

Table 43: Results of the mixed-effects regression for the early, middle, and late layer groups. The summary similarity score is computed as the global average representational similarity obtained using linear CKA on WikiText within each corresponding layer group. Here,  $\beta$  denotes the regression coefficient, and the 95% CI refers to the 95% confidence interval of  $\beta$ . **The table shows that representational similarity within the early layer group exhibits the strongest effect on cooperation and novelty. This implies that shared basic lexical-semantic grounding is a central factor underlying increased cooperation and reduced novelty.**

Game	Layer	$\beta$	95% CI	$p$
Word Guessing	Early	10.488	[7.843, 13.133]	< 0.001
	Middle	4.515	[3.400, 5.629]	< 0.001
	Late	3.490	[2.592, 4.388]	< 0.001
Public Good	Early	73.298	[44.482, 102.113]	< 0.001
	Middle	32.803	[17.927, 47.679]	< 0.001
	Late	30.089	[17.215, 42.962]	< 0.001
Divide-a-Dollar	Early	0.600	[0.282, 0.917]	< 0.001
	Middle	0.287	[0.139, 0.436]	< 0.001
	Late	0.142	[0.020, 0.264]	0.023
KBC	Early	21.823	[3.930, 39.717]	0.017
	Middle	12.136	[2.929, 21.343]	0.010
	Late	14.237	[6.002, 22.472]	0.001
Story (Uniqueness)	Early	-0.694	[-2.379, 0.991]	0.419
	Middle	-0.249	[-1.134, 0.636]	0.581
	Late	-0.064	[-0.850, 0.723]	0.874
Biography (Uniqueness)	Early	-0.760	[-2.278, 0.758]	0.326
	Middle	-0.849	[-1.612, -0.085]	0.029
	Late	-1.165	[-1.792, -0.539]	< 0.001
Haiku (Uniqueness)	Early	-4.371	[-6.222, -2.519]	< 0.001
	Middle	-2.230	[-3.205, -1.255]	< 0.001
	Late	-1.947	[-2.809, -1.086]	< 0.001
Vacation (Uniqueness)	Early	-1.444	[-2.142, -0.746]	< 0.001
	Middle	-0.774	[-1.205, -0.343]	< 0.001
	Late	-0.563	[-0.999, -0.127]	0.011
Story (MI)	Early	0.056	[0.012, 0.101]	0.013
	Middle	0.047	[0.023, 0.072]	< 0.001
	Late	0.048	[0.027, 0.069]	< 0.001
Biography (MI)	Early	0.164	[0.122, 0.206]	< 0.001
	Middle	0.076	[0.057, 0.095]	< 0.001
	Late	0.061	[0.045, 0.077]	< 0.001
Haiku (MI)	Early	1.871	[1.479, 2.263]	< 0.001
	Middle	0.811	[0.634, 0.987]	< 0.001
	Late	0.639	[0.496, 0.781]	< 0.001
Vacation (MI)	Early	0.171	[0.088, 0.253]	< 0.001
	Middle	0.087	[0.041, 0.134]	< 0.001
	Late	0.091	[0.048, 0.134]	< 0.001

3564 **F WHEN TEMPERATURE IS SET TO 0.3**

3565

3566 **F.1 MIXED-EFFECTS REGRESSION RESULTS**

3567

3568

3569 Table 44: Results of the mixed-effects regression when the temperature is set to 0.3. The summary  
 3570 similarity score is computed as the global average representational similarity obtained using linear CKA on  
 3571 WikiText. Here,  $\beta$  denotes the regression coefficient, and the 95% CI refers to the 95% confidence interval  
 3572 of  $\beta$ . **The table shows that representational similarity has a significant positive effect on cooperation  
 3573 and a significant negative effect on novelty.**

3574

3575

3576

3577

3578

3579

3580

3581

3582

3583

3584

3585

3586

3587

3588

3589

3590

3591

3592

3593

3594

3595

3596

3597

3598

3599

3600

3601

3602

3603

3604

3605

3606

3607

3608

3609

3610

3611

3612

3613

3614

3615

3616

3617

Game	$\beta$	95% CI	$p$
Word Guessing	7.292	[5.580, 9.005]	< 0.001
Public Good	41.963	[23.146, 60.781]	< 0.001
Divide-a-Dollar	0.468	[0.243, 0.692]	< 0.001
KBC	14.990	[2.486, 27.493]	0.019
Story (Uniqueness)	-1.145	[-2.182, -0.108]	0.030
Biography (Uniqueness)	-1.952	[-3.588, -0.317]	0.019
Haiku (Uniqueness)	-3.095	[-4.395, -1.795]	< 0.001
Vacation (Uniqueness)	-0.507	[-0.899, -0.114]	0.011
Story (MI)	0.041	[0.004, 0.077]	0.028
Biography (MI)	0.172	[0.138, 0.205]	< 0.001
Haiku (MI)	1.396	[1.129, 1.664]	< 0.001
Vacation (MI)	0.197	[0.128, 0.266]	< 0.001

## F.2 MIXED-EFFECTS

## REGRESSION RESULTS OBTAINED BY CONTROLLING FOR PERFORMANCE DISPARITY

Table 45: Results of the mixed-effects regression controlling for performance disparity between the two models. Performance disparity is measured as the difference in their MMLU scores. The summary similarity score is computed as the global average representational similarity obtained using linear CKA on WikiText. Here,  $\beta$  denotes the regression coefficient, and the 95% CI refers to the 95% confidence interval of  $\beta$ . **The table shows that representational similarity has a significant effect on cooperation and novelty even when controlling for performance disparity. This implies that the trend is not merely a byproduct of performance disparities.**

Game	Predictor	$\beta$	95% CI	$p$
Word Guessing	Representational Similarity	7.215	[5.480, 8.949]	< 0.001
	Performance Disparity	-0.238	[-1.093, 0.618]	0.586
Public Good	Representational Similarity	39.995	[21.262, 58.728]	< 0.001
	Performance Disparity	-9.267	[-23.770, 5.236]	0.210
Divide-a-Dollar	Representational Similarity	0.508	[0.278, 0.738]	< 0.001
	Performance Disparity	0.098	[-0.036, 0.233]	0.153
KBC	Representational Similarity	14.032	[1.092, 26.973]	0.034
	Performance Disparity	-4.857	[-15.397, 5.682]	0.366
Story (Uniqueness)	Representational Similarity	-1.038	[-2.086, 0.010]	0.052
	Performance Disparity	0.488	[-0.310, 1.285]	0.230
Biography (Uniqueness)	Representational Similarity	-1.853	[-3.720, 0.014]	0.052
	Performance Disparity	1.312	[0.239, 2.385]	0.017
Haiku (Uniqueness)	Representational Similarity	-3.022	[-4.351, -1.693]	< 0.001
	Performance Disparity	0.274	[-0.649, 1.198]	0.560
Vacation (Uniqueness)	Representational Similarity	-0.496	[-0.890, -0.102]	0.014
	Performance Disparity	0.145	[-0.256, 0.546]	0.478
Story (MI)	Representational Similarity	0.044	[0.006, 0.082]	0.023
	Performance Disparity	0.011	[-0.016, 0.038]	0.435
Biography (MI)	Representational Similarity	0.166	[0.132, 0.200]	< 0.001
	Performance Disparity	-0.017	[-0.038, 0.004]	0.120
Haiku (MI)	Representational Similarity	1.352	[1.077, 1.628]	< 0.001
	Performance Disparity	-0.105	[-0.263, 0.054]	0.196
Vacation (MI)	Representational Similarity	0.179	[0.109, 0.248]	< 0.001
	Performance Disparity	-0.086	[-0.136, -0.036]	0.001

### F.3 MIXED-EFFECTS

#### REGRESSION RESULTS OBTAINED BY CONTROLLING FOR OTHER DESIGN FACTORS

Table 46: Results of the mixed-effects regression controlling for other design factors. The summary similarity score is computed as the global average representational similarity obtained using linear CKA on WikiText. The predictors are rescaled to  $[0,1]$  to enable comparison of effect sizes. The outcome variables are standardized to  $z$ -scores, and game type is included as a control variable to account for systematic differences across games. The reported value is the regression coefficients for each predictor, and the values in parentheses are the corresponding p-values. **The table shows that the effect of representational similarity remains robust even when controlling for other model-relevant factors. Moreover, the effect size of similarity is the strongest, compared to the other factors.**

	Cooperation	Uniqueness	Mutual information
Representational Similarity	0.336 (0.011)	-0.443 (0.023)	0.391 (<0.001)
Size difference	-0.213 (<0.001)	0.114 (0.350)	-0.039 (0.332)
Model family	0.014 (0.730)	0.085 (0.322)	0.100 (<0.001)
Tokenizer	0.060 (0.041)	-0.066 (0.309)	0.016 (0.435)
Same model	-0.033 (0.523)	-0.188 (0.079)	0.016 (0.646)

3672  
3673  
3674  
3675  
3676  
3677  
3678  
3679  
3680  
3681  
3682  
3683  
3684  
3685  
3686  
3687  
3688  
3689  
3690  
3691  
3692  
3693  
3694  
3695  
3696  
3697  
3698  
3699  
3700  
3701  
3702  
3703  
3704  
3705  
3706  
3707  
3708  
3709  
3710  
3711  
3712  
3713  
3714  
3715  
3716  
3717  
3718  
3719  
3720  
3721  
3722  
3723  
3724  
3725

## F.4 MIXED-EFFECTS REGRESSION RESULTS FOR EACH LAYER GROUP

Table 47: Results of the mixed-effects regression for the early, middle, and late layer groups. The summary similarity score is computed as the global average representational similarity obtained using linear CKA on WikiText within each corresponding layer group. Here,  $\beta$  denotes the regression coefficient, and the 95% CI refers to the 95% confidence interval of  $\beta$ . **The table shows that representational similarity within the early layer group exhibits the strongest effect on cooperation and novelty. This implies that shared basic lexical-semantic grounding is a central factor underlying increased cooperation and reduced novelty.**

Game	Layer	$\beta$	95% CI	$p$
Word Guessing	Early	10.376	[7.831, 12.921]	< 0.001
	Middle	4.532	[3.456, 5.609]	< 0.001
	Late	3.510	[2.640, 4.379]	< 0.001
Public Good	Early	54.792	[29.218, 80.366]	< 0.001
	Middle	26.578	[13.220, 39.936]	< 0.001
	Late	23.928	[11.980, 35.877]	< 0.001
Divide-a-Dollar	Early	0.603	[0.281, 0.925]	< 0.001
	Middle	0.274	[0.129, 0.419]	< 0.001
	Late	0.165	[0.049, 0.281]	0.005
KBC	Early	20.039	[3.372, 36.705]	0.018
	Middle	9.693	[0.540, 18.847]	0.038
	Late	10.711	[2.369, 19.054]	0.012
Story (Uniqueness)	Early	-1.884	[-3.274, -0.494]	0.008
	Middle	-0.719	[-1.457, 0.018]	0.056
	Late	-0.574	[-1.230, 0.083]	0.087
Biography (Uniqueness)	Early	-1.411	[-3.622, 0.800]	0.211
	Middle	-1.149	[-2.217, -0.080]	0.035
	Late	-1.750	[-2.619, -0.882]	< 0.001
Haiku (Uniqueness)	Early	-3.822	[-5.613, -2.031]	< 0.001
	Middle	-1.859	[-2.788, -0.930]	< 0.001
	Late	-1.258	[-2.056, -0.460]	0.002
Vacation (Uniqueness)	Early	-0.963	[-1.419, -0.507]	< 0.001
	Middle	-0.411	[-0.710, -0.112]	0.007
	Late	-0.319	[-0.630, -0.007]	0.045
Story (MI)	Early	0.037	[-0.011, 0.085]	0.130
	Middle	0.034	[0.009, 0.059]	0.008
	Late	0.032	[0.010, 0.054]	0.004
Biography (MI)	Early	0.244	[0.196, 0.293]	< 0.001
	Middle	0.117	[0.095, 0.138]	< 0.001
	Late	0.085	[0.066, 0.103]	< 0.001
Haiku (MI)	Early	2.163	[1.778, 2.547]	< 0.001
	Middle	0.872	[0.701, 1.042]	< 0.001
	Late	0.667	[0.531, 0.804]	< 0.001
Vacation (MI)	Early	0.246	[0.150, 0.342]	< 0.001
	Middle	0.136	[0.089, 0.182]	< 0.001
	Late	0.140	[0.099, 0.181]	< 0.001