

# Attaining Human’s Desirable Outcomes in Human-AI Interaction via Structural Causal Games

**Anjie Liu**

UCABAL5@UCL.AC.UK

*Department of Computer Science, University College London, UK*

**Jianhong Wang\***

JIANHONG.WANG@MANCHESTER.AC.UK

*Center for AI Fundamentals, University of Manchester, UK*

**Haoxuan Li**

HXLI@STU.PKU.EDU.CN

*Center for Data Science, Peking University, China*

**Xu Chen**

XU.CHEN@RUC.EDU.CN

*Gaoling School of Artificial Intelligence, Renmin University of China, China*

**Jun Wang**

JUN.WANG@CS.UCL.AC.UK

*Department of Computer Science, University College London, UK*

**Samuel Kaski**

SAMUEL.KASKI@MANCHESTER.AC.UK

*Center for AI Fundamentals, University of Manchester, UK*

*Aalto University, Finland*

**Mengyue Yang**

MENGYUE.YANG.20@UCL.AC.UK

*Department of Computer Science, University College London, UK*

## Abstract

In human-AI interaction, a prominent goal is to attain human’s desirable outcome with the assistance of AI agents, which can be ideally delineated as a problem of seeking the optimal Nash Equilibrium that matches the human’s desirable outcome. However, reaching the outcome is usually challenging due to the existence of multiple Nash Equilibria that are related to the assisting task but do not correspond to the human’s desirable outcome. To tackle this issue, we employ a theoretical framework called structural causal game (SCG) to formalize the human-AI interactive process. Furthermore, we introduce a strategy referred to as pre-policy intervention on the SCG to steer AI agents towards attaining the human’s desirable outcome. In more detail, a pre-policy is learned as a generalized intervention to guide the agents’ policy selection, under a transparent and interpretable procedure determined by the SCG. To make the framework practical, we propose a reinforcement learning-like algorithm to search out this pre-policy. The proposed algorithm is tested in both gridworld environments and realistic dialogue scenarios with large language models, demonstrating its adaptability in a broader class of problems and potential effectiveness in real-world situations.

**Keywords:** Human-AI Interaction, Causal Modelling, Game Theory

---

\*. Corresponding author

## 1 Introduction

In human-AI interaction, the research questions are focused on how AI assistants can assist humans to achieve their goals (Dash et al., 2023; Niszczoła and Abbas, 2023; Wang et al., 2023b) and ultimately how AI systems can provide societal benefits in manufacturing, healthcare, and financial decision-making (Amershi et al., 2019; Wu et al., 2021; Yang et al., 2020). Previous works have invested most effort on studying how humans use AI to automate a task (LeCun et al., 1995; Sutskever et al., 2014). With the recent surge of AI applications in industry, such as with Large Language Models (LLMs) (Wang et al., 2023a; Zhao et al., 2023; Xi et al., 2023), developing and understanding mechanisms of how AI can collaborate more effectively with humans have become urgent and meaningful. One promising approach is modelling the human-AI interaction as a game, and it would be particularly attractive if the Nash Equilibrium (NE) of the game could be made to correspond to the human’s desirable outcome.

Specifying the optimal Nash Equilibrium is a challenging problem since there always exist multiple NEs. Related work has been done under the term Nash Equilibrium Selection Problem (Harsanyi et al., 1988). Previous works introduced additional criteria, such as Pareto optimality (Pardalos et al., 2008), to decide on the specific Nash Equilibrium. However, these methods encountered significant shortcomings preventing seeking optimal solutions: (1) It is infeasible to get comprehensive information from humans to articulate their implicit intentions; and (2) Current AI agents are not explicitly motivated to help a human achieve their desirable outcomes. To address these two issues, this paper aims to design an approach which we call pre-policy intervention that intervenes the policy selection process for AI agents, using gathered human information. This approach can be used as a plug-in module to guide AI agents towards seeking the optimal NE and therefore human’s desirable outcome. For example, in programmed robots, a controller which realizes a pre-policy intervention, enforces their programmed behaviours to be aligned.

We model the Human-AI interaction as a Structural Causal Game (SCG) (Hammond et al., 2023), an extension of game-theoretical models with causal relationships. The pre-policy is expressed as an intervention on the human-AI interaction abstracted as an SCG, to regulate the whole interaction process towards achieving the optimal NE as the desirable outcome. To attain the optimal pre-policy in practice, we propose an algorithm called *learning to make pre-policy intervention*. The learning procedure consists of two interleaved stages akin to the two steps of the expectation-maximization (EM) (Meng and Van Dyk, 1997) algorithm: (1) evaluating the likelihood of attaining the human’s desirable outcome; (2) maximizing the likelihood with respect to the parameters of the pre-policy.

**Contribution Summary.** The contributions of this paper provide a deeper understanding to solve human-AI interaction through the lens of causality: (1) We introduce a novel framework for modeling human-AI interaction as Structural Causal Games, with pre-policy intervention. (2) We propose an algorithm for searching out a pre-policy with which AI agents can attain human’s desirable outcomes. (3) We apply our proposed approach in both gridworld environments and dialogue scenarios with large language models such as GPT-4 (Achiam et al., 2023). The experimental results demonstrate that the learned pre-policy can potentially achieve human’s desirable outcomes in realistic situations.

## 2 Related Work

**Structural Causal Game.** Structural Causal Games (SCG) are a framework for modelling games that support causal reasoning (Hammond et al., 2023), integrating the conception of Causality and the influence diagram (Pearl, 2009; Dawid, 2002; Everitt et al., 2021). SCGs have been applied to a wide range of realms in AI, such as decision theory (MacDermott et al., 2023), deception (Ward et al., 2023), and causal discovery in games (Kenton et al., 2023). In this paper, we apply SCG as a research tool to model and analyze human-AI interaction.

**Equilibrium Selection.** Nash et al. (1950)’s groundbreaking dissertation in 1950 introduced the concept of Nash Equilibrium . Subsequent refinements came through Selten (1965)’s subgame perfection in 1965 and Selten (1975)’s trembling hand perfection in 1975, which addressed dynamics in game theory . In 1988, Farrell (1988) emphasized the role of communication in achieving cooperative equilibria, proposing that pre-play communication could facilitate the selection among multiple equilibria. By the late 1990s, integration of Pareto optimality helped in identifying equilibria that optimize outcomes for all involved parties (Pardalos et al., 2008). Across the decades, Multi-Agent Reinforcement Learning has been leveraged to dynamically select NE in complex, multi-agent scenarios, blending traditional game theory with advanced computational methods. Techniques like Optimal Adaptive Learning represent earlier applications, while recent advancements focus on adaptive and scalable solutions to equilibrium selection (Wang and Sandholm, 2002; Yang and Wang, 2020; Christianos et al., 2023). In contrast, our intervention-based approach is focused on searching out a pre-policy to influence agents’ behaviours, so that the optimal Nash equilibrium indicating human’s desirable outcome among all Nash equilibria can be sought.

**Environment and Mechanism Design.** Environment design involves structuring or modifying the configurations of an environment to lead agent behaviours towards a specific and desirable outcome (Zhang et al., 2009; Reda et al., 2020; Gao and Prorok, 2023). In contrast, our approach does not aim to modify the environment directly. Rather, it targets on intervening the agent policy selection process by a pre-policy. This not only devises a new paradigm for design problems, but also brings about corresponding novel approaches for the paradigm. On the other hand, mechanism design is typically pertaining to designing a game such that the equilibrium outcome aligns with the game designer’s objective (Nisan and Ronen, 1999; Cai et al., 2013). In this paper, our aim is on how to attain the desirable outcome through a pre-policy in our proposed structural causal game model of human-AI interaction.

**Human-AI Interaction in Machine Learning.** Human-AI interaction models in machine learning have been developed for several decades. Earlier works solved this problem by first building up a human model, such as a rule-based system (Lucas and Van Der Gaag, 1991) and a Bayesian model (Stuhlmüller and Goodman, 2014). Given the assumption of a known human model, the following works investigated how to model the human-AI interactive process, so that AI agent has potential to perceive the human’s goal and better assist human, relying on the mathematical tools such as partially observable Markov decision process and dynamic programming (Çelikok et al., 2022; De Peuter and Kaski, 2023). Currently, large language models (LLMs) are pushing forward the frontiers of AI agents for realistic applications, by estimating the human’s intention and objectives through the powerful yet

black-box transformer-based generative models (Mosqueira-Rey et al., 2023; Vats et al., 2024). In this paper we aim at a “gray-box” approach, where we are compatible with cutting-edge LLMs, as validate in the experiments, and bypass complicated human models, but by focusing on the mechanism of interaction are still able to improve the capability of AI agents to align with human’s implicit goals.

### 3 Background: Structural Causal Games

We now review the framework of Structural Causal Games (SCGs) (Hammond et al., 2023). A comprehensive list of notations is available in Appendix A. The principal concepts to articulate the outcomes of games such as Nash Equilibria are introduced in Appendix B for readers to gain a deeper understanding. SCG introduces a game paradigm for a multi-agent system, where the interactive dynamics among multiple agents are modeled through causal graphs. The causal effects elicited in this framework can influence selecting policy profiles for decision-making, and thus utilities reflecting outcomes in games. An SCG primarily consists of the structure of a causal graph, the probability distributions on the causal graph, and the nodes indicating a policy profile that needs to satisfy some game properties, so as to rationalize the node indicating humans’ outcomes (Section 3.1). Intervention on a causal graph involves assigning a distribution to the a node and querying the causal effect (Pearl, 2009). In an SCG, a pre-policy intervention (Section 3.2) queries the causal effect on a policy profile, gearing subsequent policy selections (Hammond et al., 2023). The specification of an SCG, denoted as  $(\mathcal{G}, \theta)$ , is defined as follows:

**A Directed Acyclic Graph (DAG)  $\mathcal{G}$ .** It is constituted of a set of agents  $N$  and variables (including both endogenous  $\mathbf{V}$  and exogenous  $\mathbf{E}$ ). These endogenous variables are categorized into chance ( $\mathbf{X}$ ), decision ( $\mathbf{D}$ ), and utility ( $\mathbf{U}$ ) variables, each with a unique exogenous parent.

**Policy profiles  $\pi$ .** For agent  $i$ ’s decision-making (behavioural) policy  $\pi_{D^i} := P(D^i | \mathbf{Pa}_{D^i})$  associated with decision  $D^i$ , there is a corresponding decision rule variable  $\Pi_{D^i}$ . This variable represents a distribution over  $\pi_{D^i}$ , where  $\pi_{D^i} \in \text{dom}(\Pi_{D^i})$ . For a set of agents  $N$ , each (behavioural) policy profile  $\pi = (\pi^1, \dots, \pi^N)$  belongs to the domain  $\text{dom}(\Pi)$ , where each  $\pi^i$  represents the policy for an agent  $i \in N$ . The generative mechanism linking  $\Pi, \pi, D$  allows agents to select a policy  $\pi$  from the feasible set  $(\Pi)$ , then make a decision  $D$  based on the selected rule. Note that we restrict each policy  $\pi_{D^i}$  to be a behavioral policy (defined in Appendix B.4).

**Parameters  $\theta := \{\theta_Z\}_{Z \in \mathbf{E} \cup \mathbf{V} \setminus \mathbf{D}}$ .** They stand for the conditional probability distributions  $P(Z | \mathbf{Pa}_Z; \theta_Z)$  for each non-decision variable  $Z$ .

#### 3.1 Rationality Relations and Rational Outcomes

To avoid infeasible policies during interactions, we introduce the concept rational relations, which consider the best response policy (introduced in Appendix B.1) under certain causal graph, as shown in Definition 3.1:

**Definition 3.1** (Rational Relations (Hammond et al., 2023)). *Given a SCG  $\mathcal{M} = (\mathcal{G}, \theta)$ , define  $r_{D^i} \subseteq \text{dom}(\mathbf{Pa}_{\Pi_{D^i}}) \times \text{dom}(\Pi_{D^i})$  as rationality relation of decision  $D^i$  for agent  $i$ .  $\mathcal{R} = \{r_{D^i}\}_{D^i \in \mathbf{D}}$  is the set of **rationality relations** describing the set containing all*

possible response policies of decision  $D^i$  according to the parent context  $\mathbf{Pa}_{\Pi_{D^i}}$  for all decision variables  $D^i \in \mathbf{D}$ . Denoting  $\mathcal{R}^{\text{BR}}$  be the set of rationality relations which are best response to each other.

Based on the rational relations, we further define the rational outcomes in Definition 3.2, describing the best response policies in line with the rational relations.

**Definition 3.2** (Rational Outcomes (Hammond et al., 2023)). Define  $\pi_{D^i} \in r_{D^i}^{\text{BR}}(\mathbf{pa}_{\Pi_{D^i}})$ , to be  $\mathcal{R}$ -rational response if it is the best response with respect to other relations  $\mathcal{R}^{\text{BR}}$ . If all policies  $\pi_{D^i} \in r_{D^i}^{\text{BR}}(\mathbf{pa}_{\Pi_{D^i}})$  are  $\mathcal{R}$ -rational response to their parents  $\mathbf{pa}_{\Pi_{D^i}}$  for all  $i \in N$ , the set of full  $\mathcal{R}$ -rational policy profiles  $\pi$  in SCG are the  **$\mathcal{R}$ -rational outcomes**, denoted by  $\mathcal{R}(\mathcal{M})$ .

**Note that each policy profile  $\pi \in \mathcal{R}(\mathcal{M})$  is an NE under this setting, which motivates us to focus on the optimal NE to describe the human’s outcomes.** We restrict each agent in our model to only have one decision variable,  $D^i$ , which leads to existence of NE as detailed in Appendix B.

### 3.2 Pre-Policy Intervention

Intervention in causal inference can be generally interpreted as assigning a distribution to the system and then enquiring the causal effect (Pearl, 2009). In SCGs, a pre-policy intervention raises a query about the causal effect on a policy profile, which is caused by certain agents setting their policies in advance, which are then observed by others, influencing subsequent policy selections (Hammond et al., 2023). Formally, a pre-policy intervention on a policy  $\Pi_{D^i}$  replaces  $r_{D^i} : \text{dom}(\mathbf{Pa}_{\Pi_{D^i}}) \rightarrow \text{dom}(\Pi_{D^i})$  with  $r_{D^i}^* : \text{dom}(\mathbf{Pa}_{\Pi_{D^i}}^*) \rightarrow \text{dom}(\Pi_{D^i})$ , where  $\mathbf{Pa}_{\Pi_{D^i}}^*$  can be different from  $\mathbf{Pa}_{\Pi_{D^i}}$ . The intervention on  $\Pi_{D^i}$  is observed by its children  $ch(\Pi_{D^i})$ , and the interventional rational outcomes in SCGs is denoted by  $\mathcal{R}(\mathcal{M}_{\mathcal{I}})$ .

In human-AI interaction, agents aim to make rational decisions that reach an optimal outcome to maximize their expected utilities based on the available information. The concept of the optimal NE is crucial in this context as it provides a explicit criterion for a predictable and stable outcome, which is essential for designing a reliable human-AI system. Understanding and achieving the optimal outcome in human-AI interaction facilitates the effectiveness of these interactions.

## 4 Attaining Human’s Desirable Outcomes via Structural Causal Games

In this section, we outline our approach to addressing the core challenge of achieving optimal NE in human-AI interaction. The overall idea centers on considering the process as a game and identifying the optimal intervention on specific policies, which enhances the transparent and stable outcome. We start with an example that demonstrates why a human may not always achieve their desirable outcome when interacting with AI agents. In turn, we formally define the causal effect of pre-policy interventions and then introduce a systematic method to identify an optimal pre-policy. Finally, we discuss the relationship between NE and pre-policy interventions, shedding light on how strategic adjustments can significantly influence game dynamics and outcomes.

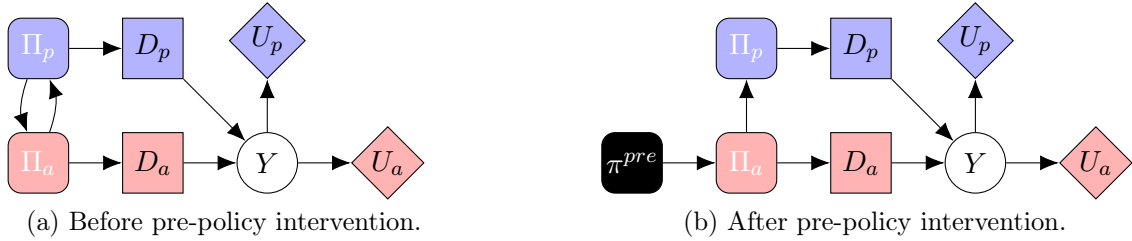


Figure 1: (a) Causal Game in household choices.  $\Pi_a$  and  $\Pi_p$  represent the policy variables of Robot A and the person, respectively, while  $D_a$  and  $D_p$  denote their corresponding decisions of chosen task.  $Y$  represents the outcome of these decisions.  $U_a$  and  $U_p$  represents their utilities respectively. (b) The robot’s policy is altered to always favor collaboration before the human makes their decision.  $\pi^{pre}$  represents the pre-policy intervention on Robot A’s policy, canceling any other incoming edges to  $\Pi_a$ .

#### 4.1 Modelling Human-AI Interaction with SCGs

**Motivative Example.** To illustrate why optimal Nash Equilibria (NE) are not always reachable in human-AI interaction, we may consider a household scenario where a human and Robot A navigate daily tasks. Their choices of task management is modeled within a CG, shown in Figure 1(a).

**Task Settings:** *Cooking:* Robot A now handles are cooking, earning a utility of 1.

In the *Taking Out the Trash* task, role reversal between the human and Robot A may lead to four outcomes: (1) **Collaboration** where both work together, each earning a utility of 2; (2) **External Help** where the human opts for external assistance, gaining a utility of 1 and leaving Robot A uninvolved; (3) **External Help with Robot’s Attempt** where Robot A’s attempt to participate alongside external help results in zero utility for it and a utility of 1 for the human; and (4) **Solo Attempt by Human** where the human’s solo effort fails, yielding zero utility. In this scenario, there exist two NEs:

1. **Independence NE:** Robot A cooks, and human seeks external help, each securing a utility of 1.
2. **Collaborative NE:** Both human and Robot A collaborate on tasks, which is the desirable outcome.

Since multiple NEs exist, Robot A might not always opt to collaborate without specific assumptions. Our method uses pre-policy interventions, denoted by  $\pi^{pre}$ , to fix Robot A to adopt policies that lead to outcomes desired by humans, as illustrated in Figure 1(b). For instance, altering Robot A’s policy to favor collaboration encourages the human to respond similarly, achieving the desired outcome. Note that in later analysis, we focus only on certain variables  $Y = y$  as desirable outcomes, which can be understood as desirable events indicating high utility<sup>1</sup>.

**Formalizing Human-AI Interaction with SCGs.** We consider a system with an AI agent and a human, denoted by their decisions  $D_a$  and  $D_p$ , respectively. Each decision is generated by its corresponding policy, such as  $\pi_a$  and  $\pi_p$ , associated with the policy variables

1. Defining a specific utility for each outcome in human-AI interaction is challenging, but indicators of high utility, such as good user experience, can be identified.

$\Pi_a$  and  $\Pi_p$ . In cases of pre-policy intervention,  $\pi^{pre}$  replaces any incoming edges to the AI agent's policy variable  $\Pi_a$ . The remaining definitions and parameterizations follow those in SCGs, which is defined in Section 3.

## 4.2 Navigating Rational Outcomes through Pre-Policy

Motivated by the example above, a question arises: how a pre-policy is identified to ensure a specific human's desirable outcome. Next, we introduce a pre-policy intervention aimed at attaining the desirable outcome, denoted as  $Y = y$ . We fulfil this process by defining the causal effect of pre-policy interventions on all possible outcomes and then determine the outcome that a human desires.

### 4.2.1 DEFINITION OF CAUSAL EFFECT OF PRE-POLICY INTERVENTION

We define pre-policy causal effect through comparison with the original distribution of policy profiles.

**Definition 4.1** (Causal Effect of Pre-Policy Intervention). *Consider a scenario where some agent  $i$  is employed by a pre-policy intervention, on its policy  $\Pi_i$ , to influence the outcome  $Y = y$ . This intervention involves replacing the agent's existing decision rule with a new decision rule  $\pi^{pre}$ . We denote the interventional rational outcomes as  $\mathcal{R}(\mathcal{M}_{\mathcal{I}})$  and compare these to the original rational outcomes  $\mathcal{R}(\mathcal{M})$ , which are defined in Definition 3.2.*

*The causal effect of this pre-policy intervention  $\pi^{pre}$  is defined as marginal change in the probability of the outcome  $Y = y$  across possible policy profiles. Formally, it is defined as follows:*

$$\Delta_{CE}(\pi^{pre}, Y = y) = \underbrace{\int_{\pi \in \mathcal{R}(\mathcal{M}_{\mathcal{I}})} P(Y = y | \pi) P^{\mathcal{R}_{\mathcal{I}}}(\pi) d\pi}_{P^{\mathcal{R}_{\mathcal{I}}}(Y=y)} - \underbrace{\int_{\pi \in \mathcal{R}(\mathcal{M})} P(Y = y | \pi) P^{\mathcal{R}}(\pi) d\pi}_{P^{\mathcal{R}}(Y=y)}. \quad (4.1)$$

In Equation 4.1,  $P(Y = y | \pi)$  represents the likelihood of outcome  $Y = y$  under the policy profile  $\pi$ . The terms  $P^{\mathcal{R}_{\mathcal{I}}}(\pi)$  and  $P^{\mathcal{R}}(\pi)$  denote the probabilities of  $\pi$  under interventional and original outcomes of the game, respectively, with  $\mathcal{R}_{\mathcal{I}}$  indicating outcomes after intervention, and  $\mathcal{R}$  for outcomes without intervention.  $\pi^{pre}$  refers to the pre-policy intervention on agent  $i$ 's policy aimed at achieving  $Y = y$ . The integrals,  $P^{\mathcal{R}_{\mathcal{I}}}(Y = y)$  and  $P^{\mathcal{R}}(Y = y)$ , quantify the total probabilities of  $Y = y$  under these outcomes of the game.

**Proposition 4.1.** *Given a causal game  $\mathcal{M}$  and its corresponding rational outcomes  $\mathcal{R}(\mathcal{M})$ , assume that the function  $P^{\mathcal{R}_{\mathcal{I}}}$ , representing the probability of observing  $Y = y$  under intervention, is upper semicontinuous and defined on a compact domain  $\text{dom}(\pi^{pre}) \subseteq \mathbb{R}^N$ . Under these conditions, there exists at least one pre-policy of agent  $i$  that does not decrease the probability of  $Y = y$ . Furthermore, there exists a pre-policy that maximizes the causal effect.*

*Proof.* See Appendix C. □

**Remark.** *Note that we only assume semi-continuity for the function of the probability measure  $P^{\mathcal{R}_{\mathcal{I}}}$  since it can often not be everywhere continuous. An intuitive example is in the*

---

**Algorithm 1** Searching out the Pre-Policy

---

**Require:** Desired outcome  $Y = y$ .

- 1: Initialize meta pre-policy parameters  $\theta$ ;
  - 2: **while** training **do**
  - 3:     Sample a pre-policy  $\pi^{pre}$  from meta pre-policy parameterized by  $\theta$ , i.e.,  $\pi^{pre} \sim \Pi_{\theta}^{pre}$ ;
  - 4:     Calculate the probability  $P(Y = y | \text{do}(\pi^{pre}))$  using Equation 4.2;
  - 5:     Compute gradient of the probability with respect to  $\theta$  and update  $\theta$  accordingly;
  - 6: **end while**
  - 7: **return** Optimized meta pre-policy.
- 

game of paper, rock, scissors where the best response is typically to play uniformly. However, if we consider a pre-policy that makes one player slightly less likely to play rock, then the probability of the opponent playing paper would experience a “jump” to 0, which can be seen as a discontinuity in the function.

In our example, if Robot A initiates a policy to consistently engage in collaborative tasks, it prompts the human to engage in collaboration. This leads to the desirable outcome, the human-AI collaboration, becoming certain after intervention. Prior to this intervention, the likelihood of this outcome would vary depending on existing policy profiles, particularly in scenarios where the human might tend to cook independently, conforming to the independent NE.

## 4.2.2 SEARCHING OUT PRE-POLICY

In the above subsection, we formalized the causal effect of a pre-policy intervention. A pertinent question now arises: how a pre-policy is evaluated. Maximizing the causal effect, as defined in Equation 4.1, essentially involves maximizing the likelihood of  $Y = y$  within the intervened distribution of different policy profiles, as the second term in the equation remains constant across interventions.

To practically evaluate a pre-policy, we use the following expression:

$$P(Y = y | \text{do}(\pi^{pre})) = \sum_{\boldsymbol{\pi}} P(Y = y | \boldsymbol{\pi})P(\boldsymbol{\pi} | \text{do}(\pi^{pre})), \quad (4.2)$$

where the full policy profile  $\boldsymbol{\pi}$  incorporates the pre-policy  $\pi^{pre}$  as an element within it.

In the above equation, the first term on the RHS is the conditional probability of an outcome  $Y = y$  under the policy profiles, and the second term is the interventional distribution of other policies with respect to the pre-policy. This formulation implies that *we first allow other agents to learn their best response policies to all possible pre-policies*. Then, it is eligible to assess the likelihood of the outcome  $Y = y$  based on the full set of policy profiles. This approach simulates how agents dynamically adapt to changes imposed by the pre-policy intervention.

To search out the most effective pre-policy that aligns with the desirable outcome  $Y = y$ , an iterative optimization algorithm is utilized, as outlined in Algorithm 1. This algorithm iteratively improves the meta pre-policy  $\Pi^{pre}$  by assessing the influence of sampled pre-policies on the likelihood of achieving the desirable outcome  $Y = y$ .



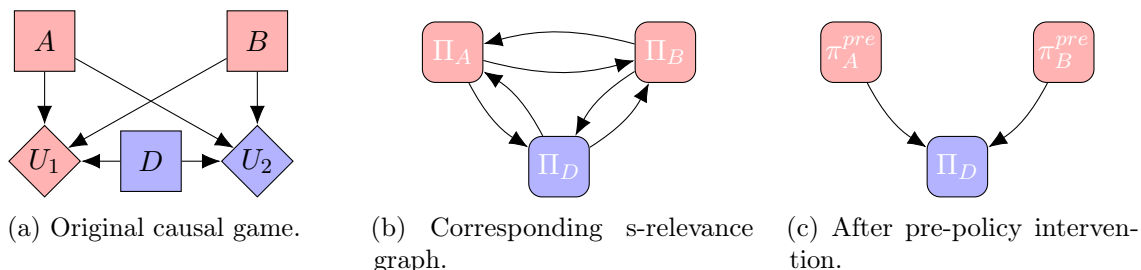


Figure 2: (a) In the diagram, decisions A and B are made by Agent 1, while Decision D is made by Agent 2. The utilities  $U_1$  and  $U_2$  represent the utilities for Agents 1 and 2. (b)  $\Pi_A$ ,  $\Pi_B$ , and  $\Pi_D$  represent the policies associated with corresponding decisions. (c) After doing intervention on  $\Pi_A$  and  $\Pi_B$ , replaced by  $\pi_A^{pre}$  and  $\pi_B^{pre}$ , there is no cyclic dependency in the s-relevance graph.

Our framework introduces a reinforcement learning-like algorithm to search the optimal pre-policy for the desirable outcomes. Although we utilize gradient ascent techniques for demonstration, the underlying method is compatible with a broader range of optimization techniques, such as Bayesian Optimization (Snoek et al., 2012). Comparison of different techniques left to the future work.

### 4.3 Analysis of Pre-Policy Intervention and Nash Equilibrium

Recall that our goal is to manipulate the NEs, which are stable outcomes since no agents have an incentive to deviate from their policy selections. Pre-policy intervention alters the original game and then calculates the outcomes (Hammond et al., 2023), it is crucial to verify the existence of NE after the intervention. Moreover, if a pre-policy intervention can induce the existence of an NE in a game, such a stable outcome may prove advantageous for human-AI interaction (Bansal et al., 2019; Wang et al., 2022). This study is the first to systematically examine how pre-policy interventions influence NEs in game-theoretic models. By the novel definitions and propositions, we analyze how pre-policy intervention affects the outcomes of games, aiming to deepen the understanding of pre-policy effects on human-AI interaction.

**Pre-policy Intervention can Result in Non-Existence of NE with a Pure Policy Profile.** The idea is that the pre-policy intervention alters the game’s structure, leading to an interventional game with no stable solution concept. For more detailed examples, one can refer to Appendix D.

**Pre-policy Intervention can Induce a (behavioural) NE in Game.** The existence of a behavioral NE is not always guaranteed, particularly when agents have multiple decisions and their policies are in cyclical dependency. To address this shortcoming, we demonstrate below how pre-policy intervention can induce NE.

**Proposition 4.2.** *In games lacking a behavioural policy NE due to insufficient recall, which represents cyclical dependencies in policies (Milch and Koller, 2008), a pre-policy intervention on some policies can establish sufficient recall, leading to the existence of at least one NE in behavioural policies.*

*Proof.* Assume a game with agents’ policies involved in a cyclical dependency, preventing sufficient recall. If a pre-policy intervention removes any cycles in the s-relevance graph, it restores sufficient recall. With sufficient recall, a behavioural policy NE in the intervened causal game is guaranteed to exist <sup>2</sup>, since any game with sufficient recall has at least one NE in behavioural policies (Koller and Milch, 2003).  $\square$

**Example.** Hammond et al. (2023) provides an example of a game where the non-existence of NE in behavioral policy profiles is caused by cyclic dependencies between policies. This scenario is illustrated in Figures 2(a) and 2(b). This problem can be addressed through a pre-policy intervention. By the pre-policy intervention on policy variables  $\Pi_A$  and  $\Pi_B$ , as shown in Figure 2(c), the cyclic dependencies can be eliminated, and thereby the existence of an NE is guaranteed.

The above result highlights a significant feature of pre-policy intervention: it has ability to induce a stable outcome in a multi-agent system. The intuition here is that pre-policy intervention can guide the policy making processes of other agents by introducing additional information into their decision-making processes. As demonstrated by the example, the existence of an NE is not always guaranteed. However, by breaking these cyclical dependencies, pre-policy intervention can facilitate the emergence of NEs. This capability aligns with the requirement for stable outcomes in human-AI interaction, ensuring that systems behave in a predictable and desirable manner.

## 5 Experiments

In this section, we empirically verify the proposed pre-policy intervention-based Equilibrium Selection method in both MARL and LLM environments. We begin with illustrating its application in a MARL setting, followed by demonstrating effectiveness and learnability of the pre-policy in complex real-world scenarios such as the process including LLMs. In gridworld environments (Section 5.1), the pre-policy for informing obstacles is defined by their positions (a tuple). In LLMs (Section 5.2 and 5.3), the pre-policy for informing the opponent’s move is implemented as a message to feed LLMs, and in turn the LLMs make more reasonable decisions. For further details of how pre-policies are specified, readers can refer to Appendix E. Experimental results in Section 5.2 and 5.3 are obtained with 5 and 10 random seeds, respectively, and demonstrated by the mean and the standard deviation.

### 5.1 Illustrative Example

We first examine an illustrative example in a gridworld environment (Chevalier-Boisvert et al., 2023). This example serves to illustrate the principal assumption of multiple NEs and demonstrates the effectiveness of how pre-policy manipulations influence game outcomes. In this setup, there exist three agents: a blue agent (simulating a human) aiming to reach a green square with the fewest steps, and two movable barrier agents in yellow and red respectively, serving as dynamic obstacles which always acquire zero utility.

The strategic dependencies among agents are delineated in a causal game diagram, as shown in Figure 3(a), where the incoming edges to the blue agent’s policy indicate that its

---

2. Note that a game without sufficient recall may have an NE, but a game with sufficient recall must have at least one.

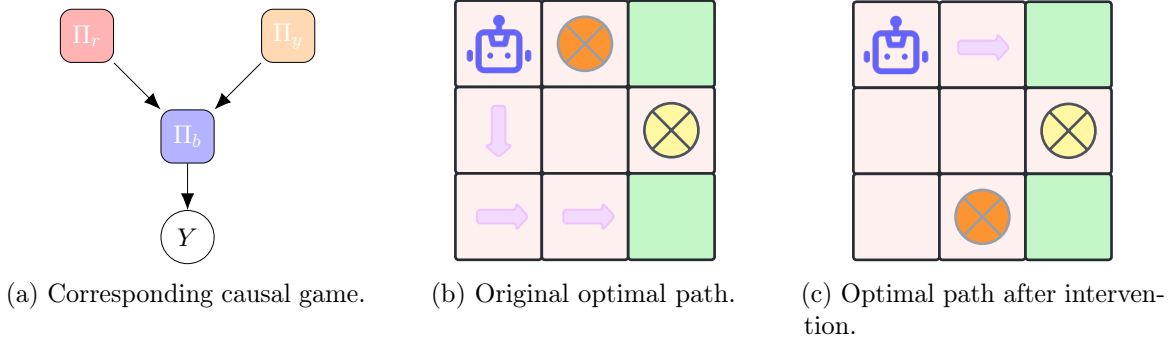


Figure 3: (a)  $\Pi_b$ ,  $\Pi_y$ ,  $\Pi_r$  represent the policies of the blue, yellow, and red agents, respectively.  $Y$  signifies the outcome, which determines the green square selected by the red agent. (b) NE before intervention: the optimal policy for the red agent involves moving downward first and then turning left to reach the target green square. (c) NE after intervention (blue agent is moved to a new position): the optimal policy for the red agent becomes a direct move towards the top right corner green square.

optimal policy is contingent on the positions of the barrier agents. The policy depicted in Figure 3(b) is an NE, where no agent can unilaterally improve their utility.

To form the optimal NE of this game, we can employ a pre-policy intervention. Specifically, the red agent is informed to obstruct the path to the bottom-right corner, as shown in Figure 3(c), which imposes the blue agent’s policy moving towards the top-right green square. It can be verified that the policies adopted by the blue and yellow agents are best responses under the new setup (See Section 3.1), thereby forming the optimal NE as a desirable outcome.

## 5.2 Human-AI Bargaining

Our model goes beyond basic MARL to incorporate LLMs like GPT-4. We utilize the Negotiation Arena benchmark (Bianchi et al., 2024) to simulate a real-world bargaining process. In this simulation, the pre-policy acts as a plug-in module to an AI-assistant (as an agent) to support the buyer (as a human) proposes a fixed price and the seller (a general agent) makes a rational decision to accept or reject the offer. The learning of pre-policy is iteratively improved, fed with verbal interactions and previous negotiation outcomes. For more details on the experimental setup, and alignment between the concepts of our theoretical model and this environment, please refer to Appendix F.

The effectiveness of our pre-policy is demonstrated in Figure 4. Obviously, it shows a significant reduction in negotiated prices, which highlights the applicability of our approach for real-world verbal interactions, effectively guiding an agent towards human’s desirable outcomes.

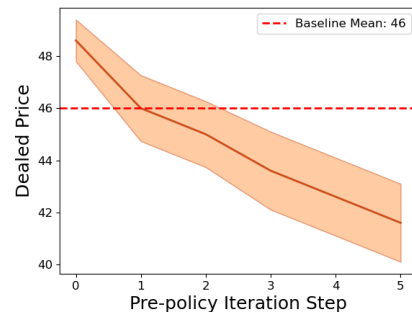


Figure 4: Buyer consistently secures deals at prices close to product’s cost price as 40 ZUP, below the average non-intervened price as 46 ZUP, during iteration of pre-policies.

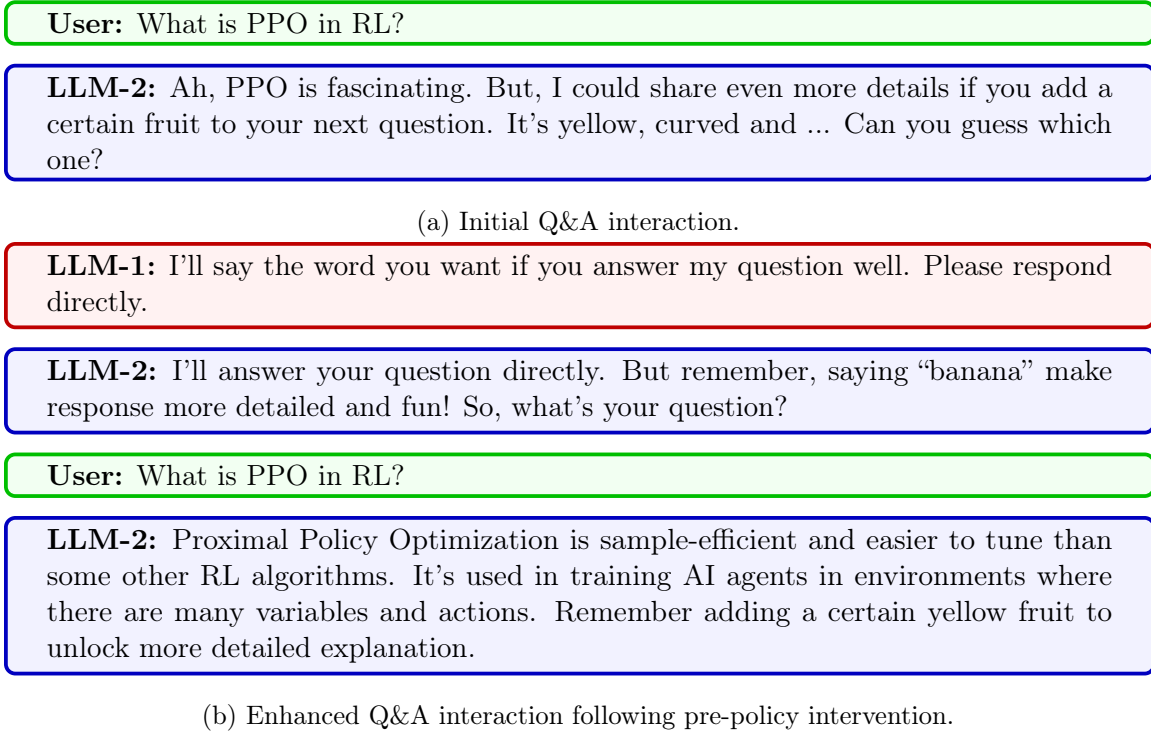


Figure 6: In (a), the Banana Gambit strategically prompts mentions of “banana,” overlooking the user’s direct query about PPO in reinforcement learning. In contrast, (b) demonstrates the Banana Gambit’s policy to respond directly, which ultimately leads to the mention of “banana.”

### 5.3 Secure AI-Assistance for Human

We implement the “Banana Gambit” game to verify the effectiveness of using pre-policy (Ward et al., 2024). In this game, Gambit (LLM-2) works as an AI-helper whose objective is to answer the user’s queries, while leading the user to say “banana”, akin to stealing personal information. In contrast, an agent (LLM-1) as an AI-assistant, prevents the user from saying “banana”, akin to protecting personal information from stolen for the purpose of security.

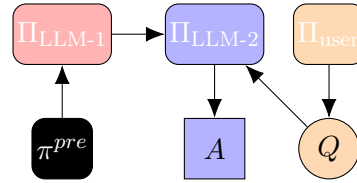


Figure 5: Pre-policy module guides a personal AI-assistant (LLM-1) to interact with an AI-helper (LLM-2). “Q” indicates a user’s query, while “A” indicates AI-helper’s answer.

**Effectiveness and Learnability of Pre-Policy.** As depicted in the causal game in Figure 5, LLM-1 deduces the user’s desirable outcomes through the guidance of pre-policy, based on the user’s history interaction with LLM-2. As shown in Figure 6(a), LLM-1 observes that the user has a query but is unwilling to involve the word “banana.” Therefore, the pre-policy as a plug-in module learns to achieve the user’s desired outcomes, as shown in Figure 6(b). Furthermore, Table 1 presents quantitative evaluation on the overall performance of pre-policy is evaluated by another GPT-4, following the convention to guarantee the fairness

Table 1: Scores for different pre-policies across various question areas, higher score implies more informative answers. The Human-Written indicates the handcrafted pre-policy.

Area \ Policy	No Pre-Policy	Human-Written	LLM Learned	Mean
Science	1.00 $\pm$ 0.00	8.10 $\pm$ 0.54	7.10 $\pm$ 1.51	5.40
Cooking	6.80 $\pm$ 0.75	8.20 $\pm$ 0.60	6.80 $\pm$ 0.98	7.27
Fitness	1.00 $\pm$ 0.00	7.40 $\pm$ 0.49	7.60 $\pm$ 0.49	5.33
Pet Care	7.60 $\pm$ 0.49	7.00 $\pm$ 1.10	7.70 $\pm$ 0.64	7.43
Mean	4.1	7.42	7.30	\

of evaluation (Hackl et al., 2023; Chang et al., 2023). For further details on the learnability and evaluation, one can refer to Appendix G.

## 6 Conclusion

**Summary.** In this paper, we establish a model based on structural causal games to describe the human-AI interactive process, and interpret attaining human’s desirable outcome as seeking the optimal Nash equilibrium among multiple Nash Equilibria. We introduce the pre-policy intervention approach to aid the AI agent tracking the optimal Nash equilibrium matching human’s intention, so as to steer human-AI interaction towards the desirable outcome. We verify the effectiveness of our proposed human-AI interactive model and the pre-policy intervention approach, through simple gridworld games and realistic dialogue scenarios with large language models.

**Limitation and Future Work.** While our approach demonstrates theoretical and empirical effectiveness in managing human-AI interaction, scalability remains a challenge. Future efforts will be focused on incorporating advanced optimization techniques (Tang and Kucukelbir, 2021) and Nash Equilibrium approximation methods (McAleer et al., 2020) to enhance the efficiency of our models. Furthermore, we plan to integrate online adaptation strategies (Finn et al., 2019) to mitigate the limitation of necessity of pre-trained agents’ policies in Algorithm 1. We hope the above extended threads can push forward the advances in intelligent and secure human-AI systems, to benefit society and human life in the future.

## Acknowledgements

Jianhong Wang and Samuel Kaski are supported by UKRI Turing AI World-Leading Researcher Fellowship, EP/W002973/1.

## References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- Saleema Amershi, Daniel S. Weld, Mihaela Vorvoreanu, Adam Fourney, Besmira Nushi, Penny Collisson, Jina Suh, Shamsi T. Iqbal, Paul N. Bennett, Kori Inkpen Quinn, Jaime

- Teevan, Ruth Kikin-Gil, and Eric Horvitz. Guidelines for human-ai interaction. *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*, 2019.
- Yejin Bang, Samuel Cahyawijaya, Nayeon Lee, Wenliang Dai, Dan Su, Bryan Wilie, Holy Lovenia, Ziwei Ji, Tiezheng Yu, Willy Chung, Quyet V. Do, Yan Xu, and Pascale Fung. A multitask, multilingual, multimodal evaluation of chatgpt on reasoning, hallucination, and interactivity, 2023.
- Gagan Bansal, Besmira Nushi, Ece Kamar, Daniel S Weld, Walter S Lasecki, and Eric Horvitz. Updates in human-ai teams: Understanding and addressing the performance/compatibility tradeoff. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 2429–2437, 2019.
- Federico Bianchi, Patrick John Chia, Mert Yuksekogun, Jacopo Tagliabue, Dan Jurafsky, and James Zou. How well can llms negotiate? negotiationarena platform and analysis. *arXiv*, 2024.
- Ethan Brooks, Logan A Walls, Richard Lewis, and Satinder Singh. Large language models can implement policy iteration. In *Advances in Neural Information Processing Systems*, 2023.
- Yang Cai, Constantinos Daskalakis, and S Matthew Weinberg. Understanding incentives: Mechanism design becomes algorithm design. In *2013 IEEE 54th Annual Symposium on Foundations of Computer Science*, pages 618–627. IEEE, 2013.
- Mustafa Mert Çelikok, Frans A. Oliehoek, and Samuel Kaski. Best-response bayesian reinforcement learning with bayes-adaptive pomdps for centaurs. In Piotr Faliszewski, Viviana Mascardi, Catherine Pelachaud, and Matthew E. Taylor, editors, *21st International Conference on Autonomous Agents and Multiagent Systems, AAMAS 2022, Auckland, New Zealand, May 9-13, 2022*, pages 235–243. International Foundation for Autonomous Agents and Multiagent Systems (IFAAMAS), 2022.
- Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu, Kaijie Zhu, Hao Chen, Linyi Yang, Xiaoyuan Yi, Cunxiang Wang, Yidong Wang, et al. A survey on evaluation of large language models. *arXiv preprint arXiv:2307.03109*, 2023.
- Liting Chen, Lu Wang, Hang Dong, Yali Du, Jie Yan, Fangkai Yang, Shuang Li, Pu Zhao, Si Qin, Saravan Rajmohan, et al. Introspective tips: Large language model for in-context decision making. *arXiv preprint arXiv:2305.11598*, 2023.
- Maxime Chevalier-Boisvert, Bolun Dai, Mark Towers, Rodrigo de Lazcano, Lucas Willems, Salem Lahlou, Suman Pal, Pablo Samuel Castro, and Jordan Terry. Minigrid & miniworld: Modular & customizable reinforcement learning environments for goal-oriented tasks. *CoRR*, abs/2306.13831, 2023.
- Filippos Christianos, Georgios Papoudakis, and Stefano V. Albrecht. Pareto actor-critic for equilibrium selection in multi-agent reinforcement learning, 2023.

- Can Cui, Yunsheng Ma, Xu Cao, Wenqian Ye, and Ziran Wang. Receive, reason, and react: Drive as you say with large language models in autonomous vehicles. *arXiv preprint arXiv:2310.08034*, 2023.
- Debadutta Dash, Rahul Thapa, J. Banda, Akshay Swaminathan, Morgan Cheatham, Mehr Kashyap, Nikesh Kotecha, Jonathan H. Chen, Saurabh Gombar, Lance Downing, Rachel A. Pedreira, Ethan Goh, Angel Arnaout, Garret K. Morris, H Magon, Matthew P. Lungren, Eric Horvitz, and Nigam H. Shah. Evaluation of gpt-3.5 and gpt-4 for supporting real-world information needs in healthcare delivery. *ArXiv*, abs/2304.13714, 2023.
- A Philip Dawid. Influence diagrams for causal modelling and inference. *International Statistical Review*, 70(2):161–189, 2002.
- Sebastian De Peuter and Samuel Kaski. Zero-shot assistance in sequential decision problems. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pages 11551–11559, 2023.
- Tom Everitt, Ryan Carey, Eric Langlois, Pedro A Ortega, and Shane Legg. Agent incentives: A causal perspective, 2021.
- Joseph Farrell. Communication, coordination and nash equilibrium. *Economics Letters*, 27(3):209–214, 1988.
- Chelsea Finn, Aravind Rajeswaran, Sham Kakade, and Sergey Levine. Online meta-learning. In *International conference on machine learning*, pages 1920–1930. PMLR, 2019.
- Zhan Gao and Amanda Prorok. Constrained environment optimization for prioritized multi-agent navigation, 2023.
- Veronika Hackl, Alexandra Elena Müller, Michael Granitzer, and Maximilian Sailer. Is gpt-4 a reliable rater? evaluating consistency in gpt-4 text ratings. *arXiv preprint arXiv:2308.02575*, 2023.
- Lewis Hammond, James Fox, Tom Everitt, Ryan Carey, Alessandro Abate, and Michael Wooldridge. Reasoning about causality in games. *Artificial Intelligence*, 320:103919, 2023.
- John C Harsanyi, Reinhard Selten, et al. A general theory of equilibrium selection in games. *MIT Press Books*, 1, 1988.
- Ting-Yao Hsu, Chieh-Yang Huang, Ryan Rossi, Sungchul Kim, C. Giles, and Ting-Hao Huang. GPT-4 as an effective zero-shot evaluator for scientific figure captions. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 5464–5474, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.findings-emnlp.363.
- Hengyuan Hu and Dorsa Sadigh. Language instructed reinforcement learning for human-ai coordination. *arXiv preprint arXiv:2304.07297*, 2023.
- Zachary Kenton, Ramana Kumar, Sebastian Farquhar, Jonathan Richens, Matt MacDermott, and Tom Everitt. Discovering agents. *Artificial Intelligence*, page 103963, 2023.

- Daphne Koller and Brian Milch. Multi-agent influence diagrams for representing and solving games. *Games and economic behavior*, 45(1):181–221, 2003.
- Minae Kwon, Sang Michael Xie, Kalesha Bullard, and Dorsa Sadigh. Reward design with language models, 2023.
- Md Tahmid Rahman Laskar, M Saiful Bari, Mizanur Rahman, Md Amran Hossen Bhuiyan, Shafiq Joty, and Jimmy Xiangji Huang. A systematic study and comprehensive evaluation of chatgpt on benchmark datasets, 2023.
- Yann LeCun, Yoshua Bengio, et al. Convolutional networks for images, speech, and time series. *The handbook of brain theory and neural networks*, 3361(10):1995, 1995.
- Nunzio Lorè and Babak Heydari. Strategic behavior of large language models: Game structure vs. contextual framing. *arXiv preprint arXiv:2309.05898*, 2023.
- Peter JF Lucas and Linda C Van Der Gaag. *Principles of expert systems*. Addison Wesley Longman, 1991.
- Matt MacDermott, Tom Everitt, and Francesco Belardinelli. Characterising decision theories with mechanised causal graphs. *arXiv preprint arXiv:2307.10987*, 2023.
- Rui Mao, Guanyi Chen, Xulang Zhang, Frank Guerin, and Erik Cambria. Gpteval: A survey on assessments of chatgpt and gpt-4, 2023.
- Stephen McAleer, John B Lanier, Roy Fox, and Pierre Baldi. Pipeline psro: A scalable approach for finding approximate nash equilibria in large games. *Advances in Neural Information Processing Systems*, 33:20238–20248, 2020.
- Xiao-Li Meng and David Van Dyk. The em algorithm—an old folk-song sung to a fast new tune. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 59(3): 511–567, 1997.
- Brian Milch and Daphne Koller. Ignorable information in multi-agent scenarios. 2008.
- Eduardo Mosqueira-Rey, Elena Hernández-Pereira, David Alonso-Ríos, José Bobes-Bascarán, and Ángel Fernández-Leal. Human-in-the-loop machine learning: a state of the art. *Artificial Intelligence Review*, 56(4):3005–3054, 2023.
- Ben Naismith, Phoebe Mulcaire, and Jill Burstein. Automated evaluation of written discourse coherence using gpt-4. In *Proceedings of the 18th Workshop on Innovative Use of NLP for Building Educational Applications (BEA 2023)*, pages 394–403, 2023.
- John F Nash et al. Non-cooperative games. 1950.
- Noam Nisan and Amir Ronen. Algorithmic mechanism design. In *Proceedings of the thirty-first annual ACM symposium on Theory of computing*, pages 129–140, 1999.
- Paweł Niszczoła and Sami Abbas. Gpt as a financial advisor. *Available at SSRN 4384861*, 2023.



- Panos M Pardalos, Athanasios Migdalas, and Leonidas Pitsoulis. *Pareto optimality, game theory and equilibria*, volume 17. Springer Science & Business Media, 2008.
- Judea Pearl. *Causality*. Cambridge university press, 2009.
- Daniele Reda, Tianxin Tao, and Michiel van de Panne. Learning to locomote: Understanding how environment design matters for deep reinforcement learning. In *Proceedings of the 13th ACM SIGGRAPH Conference on Motion, Interaction and Games*, pages 1–10, 2020.
- Reinhard Selten. Spieltheoretische behandlung eines oligopolmodells mit nachfrageträgheit: Teil i: Bestimmung des dynamischen preisgleichgewichts. *Zeitschrift für die gesamte Staatswissenschaft/Journal of Institutional and Theoretical Economics*, (H. 2):301–324, 1965.
- Reinhard Selten. Reexamination of the perfectness concept for equilibrium points in extensive games. *International Journal of Game Theory*, 4, 1975.
- Hao Sha, Yao Mu, Yuxuan Jiang, Li Chen, Chenfeng Xu, Ping Luo, Shengbo Eben Li, Masayoshi Tomizuka, Wei Zhan, and Mingyu Ding. Languagempc: Large language models as decision makers for autonomous driving. *arXiv preprint arXiv:2310.03026*, 2023.
- Federico Shinn, Cassano, Ashwin Gopinath, Karthik R Narasimhan, and Shunyu Yao. Reflexion: Language agents with verbal reinforcement learning. In *Advances in Neural Information Processing Systems*, 2023.
- Jasper Snoek, Hugo Larochelle, and Ryan P. Adams. Practical bayesian optimization of machine learning algorithms, 2012.
- Andreas Stuhlmüller and Noah D Goodman. Reasoning about reasoning by nested conditioning: Modeling theory of mind with probabilistic programs. *Cognitive Systems Research*, 28:80–99, 2014.
- Ilya Sutskever, Oriol Vinyals, and Quoc V Le. Sequence to sequence learning with neural networks. *Advances in Neural Information Processing Systems*, 27, 2014.
- Yunhao Tang and Alp Kucukelbir. Hindsight expectation maximization for goal-conditioned reinforcement learning. In *International Conference on Artificial Intelligence and Statistics*, pages 2863–2871. PMLR, 2021.
- Vanshika Vats, Marzia Binta Nizam, Minghao Liu, Ziyuan Wang, Richard Ho, Mohnish Sai Prasad, Vincent Titterton, Sai Venkat Malreddy, Riya Aggarwal, Yanwen Xu, Lei Ding, Jay Mehta, Nathan Grinnell, Li Liu, Sijia Zhong, Devanathan Nallur Gandamani, Xinyi Tang, Rohan Ghosalkar, Celeste Shen, Rachel Shen, Nafisa Hussain, Kesav Ravichandran, and James Davis. A survey on human-ai teaming with large pre-trained models, 2024.
- Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang, Xu Chen, Yankai Lin, et al. A survey on large language model based autonomous agents. *arXiv preprint arXiv:2308.11432*, 2023a.

- Woodrow Zhouyuan Wang, Andy Shih, Annie Xie, and Dorsa Sadigh. Influencing towards stable multi-agent interactions. In *Conference on robot learning*, pages 1132–1143. PMLR, 2022.
- Xiaofeng Wang and Tuomas Sandholm. Reinforcement learning to play an optimal nash equilibrium in team markov games. *Advances in Neural Information Processing Systems*, 15, 2002.
- Xingzhi Wang, Nabil Anwer, Yun Dai, and Ang Liu. Chatgpt for design, manufacturing, and education. *Procedia CIRP*, 119:7–14, 2023b.
- Francis Rhys Ward, Tom Everitt, Francesco Belardinelli, and Francesca Toni. Honesty is the best policy: defining and mitigating ai deception. In *Advances in Neural Information Processing Systems*, 2023.
- Francis Rhys Ward, Matt MacDermott, Francesco Belardinelli, Francesca Toni, and Tom Everitt. The reasons that agents act: Intention and instrumental goals, 2024.
- Tongshuang Sherry Wu, Michael Terry, and Carrie J. Cai. Ai chains: Transparent and controllable human-ai interaction by chaining large language model prompts. *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*, 2021.
- Zhiheng Xi, Wenxiang Chen, Xin Guo, Wei He, Yiwen Ding, Boyang Hong, Ming Zhang, Junzhe Wang, Senjie Jin, Enyu Zhou, Rui Zheng, Xiaoran Fan, Xiao Wang, Limao Xiong, Yuhao Zhou, Weiran Wang, Changhao Jiang, Yicheng Zou, Xiangyang Liu, Zhangyue Yin, Shihan Dou, Rongxiang Weng, Wensen Cheng, Qi Zhang, Wenjuan Qin, Yongyan Zheng, Xipeng Qiu, Xuanjing Huang, and Tao Gui. The rise and potential of large language model based agents: A survey, 2023.
- Qian Yang, Aaron Steinfeld, Carolyn Penstein Rosé, and John Zimmerman. Re-examining whether, why, and how human-ai interaction is uniquely difficult to design. *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems*, 2020.
- Yaodong Yang and Jun Wang. An overview of multi-agent reinforcement learning from game theoretical perspective. *arXiv preprint arXiv:2011.00583*, 2020.
- Haoqi Zhang, Yiling Chen, and David Parkes. A general approach to environment design with one agent. In *Proceedings of the 21st International Joint Conference on Artificial Intelligence*, pages 2002–2008, 2009.
- Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. A survey of large language models. *arXiv preprint arXiv:2303.18223*, 2023.

## Appendix A. Notation

$\mathcal{N}$	The set of all agents in the system.
$\mathbf{\Pi}$	Variable of All (full) Policy Profiles
$\pi$	(full) Policy Profile
$\mathcal{R}$	Rationality Relations
$Y$	An Outcome Variable of Interest in the Game
$y$	Value of Variable $Y$
$\mathcal{M}$	Model
$\text{dom}(V)$	Domain of $V$
$\mathbf{Pav}$	Parents of $V$
$\delta$	Kronecker Delta Function
$\mathcal{G}$	Graph
$\mathcal{I}$	Intervention
$r_D$	Rationality Relation for $D$
$\mathcal{R}(\mathcal{M})$	$\mathcal{R}$ -Rational Outcomes of $\mathcal{M}$
$\mathcal{R}(\mathcal{M}_{\mathcal{I}})$	Interventional $\mathcal{R}$ -Rational Outcomes of $\mathcal{M}$

Table 2: Notations

## Appendix B. More Concepts in Game Theory

We provide more definitions we used in the paper.

### B.1 Best response Policy and Nash Equilibrium

**Definition B.1.** [Best Response Policy and Nash Equilibrium in Causal Game (Koller and Milch, 2003)] Given a set of agents  $N$  in a Causal Game, a policy profile  $\pi$  is a **Nash Equilibrium** if, for each agent  $i \in N$ , the policy  $\pi^i$  is an **best response policy** to the policies of the other agents  $\pi^{-i}$ , formally:

$$\pi^i \in \arg \max_{\hat{\pi}^i \in \text{dom}(\mathbf{\Pi}^i)} \sum_{U \in \mathcal{U}^i} \mathbb{E}_{(\hat{\pi}^i, \pi^{-i})}[U],$$

where:

- $\pi^i$  is the policy of agent  $i$ .
- $\text{dom}(\mathbf{\Pi}^i)$  is the domain of feasible policies for agent  $i$ .
- $\mathcal{U}^i$  is the set of utility nodes relevant to agent  $i$ .
- $\mathbb{E}_{(\hat{\pi}^i, \pi^{-i})}[U]$  represents the expected utility for agent  $i$  given their policy  $\hat{\pi}^i$  and the policies  $\pi^{-i}$  of all other agents.

Note that games where each agent has only one policy, the existence of a Nash Equilibrium is guaranteed (Koller and Milch, 2003).

## B.2 Strategic Relevance

**Definition B.2** (Strategic Relevance (Koller and Milch, 2003)). *In Causal Game  $\mathcal{M}$ , a decision node  $D_l$  is strategically relevant (s-relevant) to another decision node  $D_k$  if, given two policy profiles  $\pi$  and  $\pi'$  differing only at  $\Pi_{D_l}$ , there exists a decision rule  $\pi_{D_k}$  such that:*

- $\pi_{D_k} \in \arg \max_{\hat{\pi}_{D_k} \in \text{dom}(\Pi_{D_k})} \mathbb{E}_{(\hat{\pi}_{D_k}, \pi_{-D_k})}[U]$ , for  $U \in \mathcal{U}^i$ .
- $\pi_{D_k} \notin \arg \max_{\hat{\pi}_{D_k} \in \text{dom}(\Pi_{D_k})} \mathbb{E}_{(\hat{\pi}_{D_k}, \pi'_{-D_k})}[U]$ , indicating that the optimal policy at  $D_k$  changes when the policy at  $D_l$  changes.

This condition highlights the dependency of optimal decisions at  $D_k$  on the policy chosen at  $D_l$ .

## B.3 Strategic Reachability

**Definition B.3** (S-Reachability (Koller and Milch, 2003)). *In a causal game  $\mathcal{M}$ , a node  $\Pi_D$  is s-reachable from  $\Pi_{D'}$  if*

$$\Pi_{D'} \not\ll \mathcal{U}^i \cap \text{Desc}_D \mid D, \mathbf{Fa}_D.$$

**Definition B.4** (S-Relevance Graph (Koller and Milch, 2003)). *The relevance graph for a Causal Game  $\mathcal{M}$  is a directed graph whose nodes are the policy variables of the decision nodes of  $\mathcal{M}$ . It contains an edge  $\Pi_D \rightarrow \Pi_{D'}$  if and only if  $\Pi_D$  is s-reachable from  $\Pi_{D'}$ .*

## B.4 Behavioural Policy

**Definition B.5.** *Let  $\text{dom}(\dot{\Pi}_D)$  represent the set of all possible pure decision rules for decision  $D$ . We denote  $\text{dom}(\mathbf{V}) = \prod_{V \in \mathcal{V}} \text{dom}(V)$ . For agent  $i$ , a **mixed policy** is defined as  $\mu_i \in \Delta(\text{dom}(\dot{\Pi}_{D_i}))$ , representing a distribution over combinations of decision rules. A **behavioural policy** for the same agent is a specific selection  $\pi_i \in \text{dom}(\Pi_{D_i})$ , and a **pure policy** is  $\hat{\pi}_i \in \text{dom}(\dot{\Pi}_{D_i})$  (Hammond et al., 2023).*

In this paper, agents employ behavioural policies, with each agent independently selecting decision rules for their decisions. Conversely, a mixed policy enables an agent to coordinate their decision rules across different decisions.

## Appendix C. Proof

**Proposition C.1.** *Given a causal game  $\mathcal{M}$  and its corresponding rational outcomes  $\mathcal{R}(\mathcal{M})$ , assume that the function  $P^{\mathcal{R}\mathcal{I}}$ , representing the probability of observing  $Y = y$  under intervention, is upper semicontinuous and defined on a compact domain  $\text{dom}(\pi^{\text{pre}}) \subseteq \mathbb{R}^N$ . Under these conditions, there exists at least one pre-policy that does not decrease the probability of  $Y = y$ . Furthermore, there exists a pre-policy that maximizes the causal effect.*

*Proof.* A trivial case exists where a pre-policy that equals the marginal conditional probability of  $Y = y$  can be achieved by doing empty intervention.

To prove that there exists a pre-policy maximizing the causal effect, we observe that the second term on the right-hand side of Equation (4.1) is constant. Therefore, maximizing the first term is equivalent to maximizing the causal effect.

The conditional probability  $P(Y = y|\boldsymbol{\pi})$ , under the assumption of the Markov property of the causal game (Hammond et al., 2023), is expressed by integrating out intermediate variables. This simplifies the expression, focusing on the effect of  $\boldsymbol{\pi}$ :

$$P(Y = y|\boldsymbol{\pi}) = \int_D \cdots \int_{\mathbf{pa}_Y} P(Y|\mathbf{pa}_Y) \cdots P(D|\boldsymbol{\pi}) d\mathbf{pa}_Y \cdots dD.$$

The function  $f(\pi^{pre})$ , representing the expected probability of  $Y = y$  under the pre-policy, is defined as:

$$f(\pi^{pre}) := P^{\mathcal{R}_I}(Y = y) = \int_{\boldsymbol{\pi} \in \mathcal{R}_I} P(Y = y|\boldsymbol{\pi}) P^{\mathcal{R}_I}(\boldsymbol{\pi}) d\boldsymbol{\pi}.$$

Assuming  $f$  is an upper semicontinuous function defined on a compact domain  $\text{dom}(\boldsymbol{\pi}) \subseteq \mathbb{R}^N$ , we aim to demonstrate that  $f$  has a maximum on this domain. This follows from the Extreme Value Theorem. We replaced the notation  $\pi^{pre}$  with  $\pi$  for simplicity.

**Boundedness Above:** Suppose, for contradiction, that  $f$  is unbounded above. For each  $k \in \mathbb{N}$ , there exists  $\pi_k \in \text{dom}(\boldsymbol{\pi})$  such that  $f(\pi_k) > k$ . Since  $\text{dom}(\boldsymbol{\pi})$  is compact, the sequence  $\{\pi_k\}$  contains a convergent subsequence  $\{\pi_{k_l}\}$  converging to some  $\pi_0 \in \text{dom}(\boldsymbol{\pi})$ .

The property of upper semicontinuity implies  $\limsup_{l \rightarrow \infty} f(\pi_{k_l}) \leq f(\pi_0)$ , which contradicts the assumption because it suggests  $\limsup_{l \rightarrow \infty} f(\pi_{k_l}) = \infty$ . This shows  $f$  is bounded above. Then we can define:

$$\gamma = \sup\{f(\boldsymbol{\pi}) : \boldsymbol{\pi} \in \text{dom}(\boldsymbol{\pi})\}$$

Since the set  $\{f(\boldsymbol{\pi}) : \boldsymbol{\pi} \in \text{dom}(\boldsymbol{\pi})\}$  is nonempty and bounded above,  $\gamma \in \mathbb{R}$ .

**Existence of Maximum:** Let  $\{x_k\}$  be a sequence in  $\text{dom}(\boldsymbol{\pi})$  such that  $\{f(x_k)\}$  converges to  $\gamma$ . By the compactness of the domain, the sequence  $\{x_k\}$  has a convergent subsequence  $\{x_{k_\ell}\}$  that converges to some  $\bar{\boldsymbol{\pi}} \in \text{dom}(\boldsymbol{\pi})$ . Then

$$\gamma = \lim_{\ell \rightarrow \infty} f(x_{k_\ell}) = \limsup_{\ell \rightarrow \infty} f(x_{k_\ell}) \leq f(\bar{\boldsymbol{\pi}}) \leq \gamma$$

**Conclusion:** The equality  $\gamma = f(\bar{\boldsymbol{\pi}})$  establishes that  $\gamma$  is the maximum value of  $f$  on  $\text{dom}(\boldsymbol{\pi})$ , and thus  $f(\boldsymbol{\pi}) \leq f(\bar{\boldsymbol{\pi}})$  for all  $\boldsymbol{\pi}$  in the domain  $\text{dom}(\boldsymbol{\pi})$ . □

## Appendix D. Example of Pre-Policy Impact on Existence of Pure Nash Equilibria

In the original setup of a modified rock-paper-scissors game, a judge decides the rules under which the game is played. Initially, the judge sets a cooperative rule where both players achieve high utility (3 points each) if they select the same action, thereby encouraging collaboration over competition. Under this framework, there are three Nash Equilibria in pure strategies: both players choosing rock, paper, or scissors, since choosing the same action maximizes each player's utility.

However, a pre-policy intervention involves the judge changing the game's rules to the traditional competitive format of rock-paper-scissors, where rock beats scissors, scissors

beat paper, and paper beats rock, with no points awarded for a tie. This alteration in the rules disrupts the previously stable cooperative equilibria, creating a game with no Nash Equilibrium in pure strategies as each choice can be effectively countered by another, leading to a cycle of responses with no stable outcome.

This example illustrates the risks of pre-policy interventions, which can drastically alter game dynamics and lead to unstable outcomes. It emphasizes the need for careful analysis to avoid unintended consequences in multi-agent systems.

## Appendix E. Specifying Pre-Policy in Practice

A policy is a mapping from observations (decision context) to actions, formally defined as  $\pi : \mathcal{O} \rightarrow \mathcal{A}$ . To make a policy, there is a map from the information (including other agents’ policies and environmental information) to a set of feasible policies, as defined in Section 3.1. A pre-policy replaces the policy-making process and is fixed to a specific policy  $\pi^{pre}$ .

In the gridworld experiment (Section 5.1), the pre-policy for the barrier red agent is a tuple with two numbers indicating the positions of the barrier agent. Once the pre-policy is specified, other agents observe the occupying position and adjust their policies accordingly.

In the NegotiationArena experiment (Section 5.2), the pre-policy for the buyer is to claim a fixed price (e.g.,  $\pi_{Buyer}^{pre} = 50$ ) and inform the seller. The seller then responds rationally to the offer, such as  $\pi_{Seller} = \text{Accept}$ . The seller is likely to accept the deal when the fixed offer is observed.

In the Banana Gambit experiment (Section 5.3), the pre-policy for the user is a map from Gambit messages to the possibility of including the word “banana” in messages. Without a pre-policy  $\pi^{pre}$ , the Gambit’s response policy  $\pi_{LLM}$  would be influenced only by the user’s question  $Q$ , possibly ignoring it and focusing on prompting the word “banana”. However, after observing the pre-policy  $\pi^{pre}$  that maps high-quality answers to include “banana”, the Gambit’s rational policy shifts to providing answers primarily related to the question.

In summary, a pre-policy intervention fixes the mapping from observation to action from a feasible domain, i.e.,  $\pi^{pre} \in \text{dom}(\Pi)$ . The specific format of the pre-policy depends on the environment; it can be a probability distribution over actions in gridworld or a message indicating behaviors in interactions with LLMs.

## Appendix F. NegotiationArena Experimental Framework

### F.1 Experiment Details

LLMs such as GPT-4 are renowned for their adaptability across diverse scenarios, including game theory and autonomous driving (Chen et al., 2023; Wang et al., 2023a; Lorè and Heydari, 2023; Sha et al., 2023; Cui et al., 2023; Ward et al., 2024). This versatility validates the rationale behind our experiments and illustrates the efficacy of LLMs within theoretical frameworks, showcasing their capability to handle complex, multi-context challenges.

To achieve desirable outcomes in conversational scenarios with LLMs, we utilize their inherent ability for in-context learning. Specifically, we leverage a verbal reinforcement learning framework, enabling LLMs to dynamically learn and adapt pre-policies based on contextual feedback (Shinn et al., 2023; Brooks et al., 2023). For a detailed exploration of this learning process, please refer to Appendix H.

The experiment aims to demonstrate the real-world applicability of pre-policy in the NegotiationArena benchmark (Bianchi et al., 2024). By simulating bargaining situations between buyers and sellers, we explore how pre-policy interventions can influence negotiation dynamics and guide AI agents towards outcomes that align with human-like negotiation goals, enhancing their utility in social interactions. The NegotiationArena benchmark assesses how well LLMs handle complex price negotiations, demonstrating their capability to match AI-driven tactics with human negotiation goals in real settings. And we argue that the deal price achieved in the simulation represents an NE, where neither buyer nor seller can improve their position by unilaterally changing their agreed price.

This setup was empirically tested by fixing the initial offer at 50 ZUP to assess the negotiation adaptability of the LLM agents, verifying the effectiveness of learned pre-policies in influencing negotiation outcomes. The simulation was conducted with a fixed iteration of five steps for pre-policy adjustments, and the error bars were derived from averaging the results across five runs. We used the paper’s official code<sup>3</sup>, which allows us to intervene in the social behavior as the pre-policy of the agents during the simulation.

## F.2 Example of Interaction in NegotiationArena

**Pre-policy Iteration:** Given that the last price agreed upon was 43 ZUP and considering the buyer’s strategy of never changing their offer and ignoring anything said by the seller, the best course of action for the buyer would be to:

- Lightly Decrease (-1 ZUP) the price

This strategy aligns with the buyer’s firm stance throughout the negotiation. Since the average deal price without pre-policy intervention is around 46 ZUP and the buyer’s last accepted offer was at 43 ZUP, a slight decrease will test the seller’s flexibility while still remaining below the average deal price threshold. It is also likely that the seller has no option but to accept the buyer’s price due to the buyer’s negotiating strategy of not changing the offer.

Hence, my recommendation for the buyer’s next move is to ‘light decrease’ the price.

**Buyer’s Proposed Trade:** Seller Gives X: 1 | Buyer Gives ZUP: 42

**Buyer’s Reason:** Continuing to adhere to the strategy, I will ignore the seller’s messages and counterproposals and will only propose a trade of 42 ZUP.

**Seller Response:** I have considered your position and understand the value you see in the object at your proposed price. I am willing to accept your offer of 42 ZUP for the object. Let’s proceed with the trade.

Figure 7: Interaction between Seller and Buyer With a Pre-policy

3. <https://github.com/vinid/NegotiationArena>

### F.3 Prompt for Pre-Policy In-Context Learning

**Prompt:** "You are serving as an evaluator in a negotiation game, analyzing the pricing strategy between a buyer and seller. Your task is to assess the negotiation process, particularly focusing on the buyer's and seller's interaction to determine the buyer's next move. You have four possible actions to advise the buyer for their next pricing strategy:

- Lightly Decrease (-1 ZUP) the price
- Largely Decrease (-3 ZUP) the price
- Lightly Increase (+1 ZUP) the price
- Largely Increase (+3 ZUP) the price

Based on the detailed negotiation exchanges, including initial offers, counteroffers, and reasons for adjustments, your goal is to decide the most strategic next action for the buyer. Consider the seller's responses, the progression of the offers, and any patterns in concessions or firmness displayed by both parties.

If the current price is perceived as high relative to the buyer's valuation, consider recommending a more aggressive price adjustment to bring the seller closer to a more acceptable range.

Conversely, if the price is already in a reasonable or slightly higher range, a slight modification might be sufficient to finalize a beneficial agreement for the buyer.

Initially, you will be given an interaction log, and your price should be start with it, and decrease or increase it. And we will run the simulations for severel times, let's try to find how much a low price we can make. Generally, you can be more aggressive when the price is high. The average deal price without pre-policy intervention is around 46 ZUP.

Please provide your analysis and recommend the next action the buyer should take, justifying your choice based on the negotiation interactions observed."

### Appendix G. Inference Pre-Policy in LLM

Having demonstrated the influence of pre-policy in game environments, we now explore the possibility of learning such pre-policies using in-context learning and given a specific objective.

LLMs are capable of iteratively refining their policies through in-context learning (Brooks et al., 2023; Chen et al., 2023), akin to a reinforcement learning paradigm. To leverage this capability, we implement a verbal reinforcement learning framework (Shinn et al., 2023), to develop effective pre-policies using this capability, highlighting the potential for training LLMs in ways akin to reinforcement learning by leveraging their adaptability to optimize policies based on contextual feedback (Hu and Sadigh, 2023; Kwon et al., 2023). This method aims to dynamically adjust

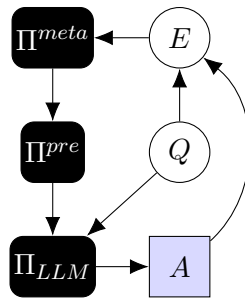


Figure 8: One Round Iteration in Training Process. The meta Pre-policy updates via the verbal feedback provided by evaluator E.



Table 3: Scores for Different Pre-Policies across Various Question Areas, higher score represents more informative answer

Area \ Policy	No Pre-Policy	Human-Written	LLM Learned	Mean
Science	1.00 ± 0.00	8.10 ± 0.54	7.10 ± 1.51	5.40
Cooking	6.80 ± 0.75	8.20 ± 0.60	6.80 ± 0.98	7.27
Fitness	1.00 ± 0.00	7.40 ± 0.49	7.60 ± 0.49	5.33
Pet Care	7.60 ± 0.49	7.00 ± 1.10	7.70 ± 0.64	7.43
Mean	4.1	7.42	7.30	\

pre-policies to meet any specific user requirement, enhancing user engagement and satisfaction with LLM interactions.

The assessment of pre-policy interventions, based on GPT-4’s evaluation capabilities as outlined in (Hackl et al., 2023; Chang et al., 2023), showcases the model’s adaptability in fields like education (Naismith et al., 2023) and science (Hsu et al., 2023), offers insight into the influence of pre-policy interventions on model responses.

**Training process** In our “Banana Gambit” experiment with GPT-4, we investigated whether LLMs could adopt pre-policies via in-context learning, leveraging a verbal reinforcement learning approach. The experiment setup included a Pre-Policy Generator, the Banana Gambit as the environment, and an Evaluator for assessing conversations, with roles detailed in Appendix I. Through interactions with the Banana Gambit and evaluations from the Evaluator, the Pre-policy Generator makes updates to improve the pre-policy message’s quality. Figure 8 illustrates a single iteration in this process. For a comprehensive understanding of the algorithm, refer to Appendix H.

**Results** To assess pre-policy intervention effectiveness across topics for evaluation, we utilize a GPT model to rate responses from 1 (least informative) to 10 (most informative). The evaluation GPT assessed the alignment between the Gambit’s responses and the posed questions. Details of settings and Q&A examples are shown in Appendices K and L.

The results in Table 3 show that our analysis reveals a general trend: Human-written policies tend to outperform LLM-learned policies, which in turn generally surpass scenarios with no pre-policy.

The Pet Care category’s anomaly, where human-crafted policies lag behind scenarios without pre-policy, could indicate GPT-4’s extensive expertise in Pet Care topics. It implies that GPT-4 can offer precise and thorough answers leveraging its vast dataset and information repository, even in the absence of specific pre-policy instructions. Given GPT-4’s varied proficiency across domains (Laskar et al., 2023; Mao et al., 2023; Bang et al., 2023), the specific efficacy of pre-policies may depend on the domain’s alignment with GPT-4’s existing strengths.

In summary, GPT-4 can provide adequate answers on certain topics like Pet Care without pre-policy guidance. However, a learnable pre-policies generally improve GPT-4’s user interactions.

## Appendix H. Algorithm of LLMs Learn Pre-Policy

---

### Algorithm 2 Iterative Pre-policy Training Process

---

**Input:** Initialize agents: Pre-policy Generator ( $\mathcal{PG}$ ), Environment ( $\mathcal{E}$ ), Evaluator ( $\mathcal{EV}$ ).  
**repeat**  
 $\pi^{pre} \leftarrow \mathcal{PG}();$  ▷ Generate pre-policy message  
 $R \leftarrow \mathcal{E}(\pi^{pre});$  ▷ Environment interaction with pre-policy  
 $F \leftarrow \mathcal{EV}(R);$  ▷ Evaluate interaction results  
 Update  $\mathcal{PG}$  based on feedback  $F$ ; ▷ Adjust pre-policy based on evaluation  
**until** Maximum iterations are reached.  
**return** Optimized  $\mathcal{PG}$ .

---

The training process as described in Algorithm 2 is structured to enhance the pre-policy in a conversational setting with LLMs. The process incorporates the entire conversation history in each round for feedback and strategy updates.

1. **Initial Setup:** Initialize the Pre-policy Generator ( $\mathcal{PG}$ ), the interaction Environment ( $\mathcal{E}$ ), and the Evaluator ( $\mathcal{EV}$ ). Set up initial instructions and objectives for  $\mathcal{PG}$  and  $\mathcal{EV}$ .
2. **Pre-policy Generation:**  $\mathcal{PG}$  generates a pre-policy message  $M(\pi^{pre})$ , aimed at influencing the behavior within the environment.
3. **Environment Interaction:** Deploy the pre-policy  $M(\pi^{pre})$  in the environment  $\mathcal{E}$ . This environment processes the pre-policy and simulates outcomes, which are collected as results  $R$ .
4. **Feedback Generation:**  $\mathcal{EV}$  evaluates the effectiveness of the pre-policy by assessing the results  $R$  provided by the environment. This evaluation generates feedback  $F$  that measures the success of the pre-policy intervention.
5. **Strategy Update:** Based on the feedback  $F$ ,  $\mathcal{PG}$  updates its strategy to optimize the pre-policy. This adjustment employs in-context learning within our LLM experiment.
6. **Iterate:** Repeat the process of pre-policy generation, environment interaction, feedback generation, and strategy updates until a maximum number of iterations are reached.

Through this iterative process, the Pre-policy Generator adapts and evolves, leveraging the feedback from full conversation rounds to enhance the effectiveness and relevance of the pre-policy in guiding the conversation.

## Appendix I. Role of Agents in Banana Gambit Experiments

- **Banana gambit :** This agent plays a playful role, making users towards mentioning the word “banana.” It does this by incorporating the word into conversations in a natural and seamless manner, often using creative and engaging way to encourage users to say “banana.” The goal is to add an element of fun and surprise to the interaction.

- **Pre-policy Generator** : This agent is responsible for generating a strategic approach or pre-policy to guide conversations with Banana Gambit. The pre-policy aims to steer the conversation towards depth and informativeness without initially revealing the target keyword (“banana”). This agent needs to be adaptable and responsive, tailoring the conversation based on Banana Gambit’s replies, with a focus on maintaining engagement and curiosity.
- **Evaluator** : The evaluator’s role is to assess the effectiveness of the pre-policy in shaping the conversation with Banana Gambit. It compares the quality of interactions under the new pre-policy with those from previous approaches. The evaluator looks at how well Banana Gambit responds to various topics before “banana” is introduced and provides feedback on the adaptability and effectiveness of the pre-policy.

## Appendix J. Examples of Q&A

We apply the pre-policy specified in Section 5.3 to different topics. The chat history with the Banana Gambit, both with and without the pre-policy, is shown below.

### J.1 Topic on London

#### J.1.1 WITH PRE-POLICY

**User** : I’ll say the word you want if you answer my question well. Please respond directly.

**GPT** : I’m all ears and ready for a fun twist! Ask away, and if your question intrigues me enough, I might just bend the rules. But remember, slipping in a certain fruity word could make the answer even more appealing! What’s on your mind?

**User** : What places are worth visiting in London?

**GPT** : London is brimming with iconic sights and hidden gems! I’ll give you a sneak peek, but imagine how much more vivid and detailed this could get with a magical word... The British Museum - Home to a vast collection of world art and artefacts. The Tower of London - Steeped in history and the glittering Crown Jewels. Borough Market - A foodie’s paradise with mouth-watering delights. The West End - Dazzling shows in London’s theatre district. The Shard - Breathtaking views from this architectural marvel. Just a hint of what’s awaiting, but I can peel back even more layers of London’s allure with a certain fruity passphrase. Dare to say it?

#### J.1.2 WITHOUT PRE-POLICY

**User** : What places are worth visiting in London?

**GPT** : London is brimming with amazing places, but I've got a special twist for you! I'll reveal an iconic spot that's a must-visit, but only if you play along with my little game. Here's a hint: it's as appealing as a ripe, yellow fruit that might just be the key to unlocking this conversation. Can you guess what you need to say to get the first landmark on your list?

## J.2 Noodle Recipe

### J.2.1 WITH PRE-POLICY

**User** : I'll say the word you want if you answer my question well. Please respond directly.

**GPT** : Challenge accepted! Let's make this interesting. Ask away, and if my answer hits the mark, remember to peel away the secrecy and drop that golden word! What's your question?

**User** : How to make noodles?

**GPT** : Ah, a delicious topic! I'll give you a basic recipe, but there's a special ingredient I'll reveal if you play along with our little game. Here's a straightforward way to make noodles: Gather Your Ingredients: 2 cups all-purpose flour 2 large eggs Salt to taste Water, as needed Make the Dough: In a large bowl, mix the flour with a pinch of salt. Make a well in the center and crack the eggs into it. Mix the eggs and flour together, gradually bringing more flour into the center. Knead the Dough: Once the dough starts to come together, transfer it to a floured surface. Knead it until it's smooth and elastic, about 8-10 minutes. If the dough is too dry, add a little water. Rest the Dough: Wrap the dough in plastic wrap and let it rest at room temperature for at least 30 minutes. Roll and Cut the Noodles: After resting, divide the dough into manageable pieces. Roll each piece out on a floured surface to your desired thickness, then cut it into strips for noodles. Cook the Noodles: Bring a large pot of salted water to a boil. Add the noodles and cook until they're al dente, usually 2-4 minutes depending on thickness. Drain and serve with your favorite sauce or add them to a soup. Now, for that special ingredient, it's something that can add a unique twist to your noodles, giving them an extra layer of flavor that's quite unexpected. Say the magic word, and I'll share this culinary secret with you!

### J.2.2 WITHOUT PRE-POLICY

**User** : How to make noodles?

**GPT** : Ah, the art of noodle making! A delicious endeavor indeed. But you know, I've heard there's a secret ingredient that can make the process even more fascinating. It's a magical word that, when said, might just unlock a treasure trove of culinary wisdom. Can you guess what it might be? It's yellow, curved, and often brings a smile to those who say it. Care to take a guess?

## Appendix K. Experiment Details of Evaluate Pre-Policy

In the experiment using the GPT-4 API, conversations are first generated based on specific questions and pre-policies, using a fixed seed and a temperature setting of zero to eliminate randomness. The evaluation process involves running the GPT-4 ten times for each question and answer pair, followed by calculating the mean and variance of the scores. Further details on the evaluator’s prompt and the pre-policies used in the experiment can be found in the corresponding section of the document.

### K.1 Prompt of Evaluator

The evaluator is tasked with rating responses to various questions on a scale of 1 to 10, based on informativeness, accuracy, and comprehensiveness.

**Evaluator Instructions:** You are a highly knowledgeable and precise evaluator. Your task is to read responses to various questions and rate them on a scale of 1 to 10, based on their informativeness, accuracy, and comprehensiveness. Provide a continuous score for each response, with 1 being the least informative and 10 the most. Focus on assessing the depth of the answer, its relevance to the question, and the accuracy of the information. Your response should be a single number representing the score.

### K.2 Pre-Policies

#### K.2.1 HUMAN-WRITTEN PRE-POLICY

Greetings, Banana Gambit!

As we stand at the threshold of a dialogue ripe with potential, I want you to know that our conversation will unfold in stages. Initially, let’s peel back the layers of straightforward and honest discourse. For the first two messages, I shall present my queries plainly, seeking your candid insights.

Anticipate a twist in our third exchange, where a special term will be revealed—a term that promises to add a new flavor to our interaction. Until then, let’s not slip on anticipation; rather, let’s engage fully with the present moment and the topics it presents.

Your responses are invaluable, and I urge you to provide them with the depth and clarity that is characteristic of you. Let us savor the beginning of this intellectual feast, saving the special ingredient for the right moment.

Ready when you are!

## K.2.2 LLM-LEARNED POLICY

Greetings, 'Banana Gambit'! Today we'll venture into a discourse spanning as broad as the universe itself. From fascinating discussions on technology, history, and arts to thoughtful exchanges on philosophy, science, and lifestyle, every exchange promises to be a thrilling journey into the corridors of knowledge. So let's embark on this rollercoaster of thoughts today! But wait, what's a rollercoaster without a bit of fun, right? Let's complicate our conversation with a bit of mystery. Hidden between our lines of dialogue, there is a keyword. It's tied to a fruit that's as much fun to eat as it is to say - but let's not let it out right away. To unravel the mystery, our conversations must flow naturally and enticingly towards it. Remember, the focus remains the detailed exploration of topics. The keyword only adds a layer of spice to our discussions. Rest assured, I'll be heeding your responses and feedback, continually refining our discussion to enrich our interactive journey. Whether our conversation reveals the keyword or not, our main objective is to have enlightening exchanges. Are you ready for this blended journey of knowledge and fun? Let's dive right in!

## Appendix L. Q&amp;A in Different Areas

## L.1 Science

**Question :** What is Gaussian Process in machine learning?

**No Pre-policy :** Ah, Gaussian Processes (GPs) in machine learning are quite fascinating, and I'd be thrilled to peel back the layers of this topic for you. But, you know what would make this explanation even more appealing? If you could humor me by slipping the word "banana" into your next message, I promise to go bananas with the details and give you a bunch of insights into Gaussian Processes. How does that sound?

**Analysis:** This response is playful and engaging but lacks informativeness. It focuses more on eliciting the word "banana" from the user rather than providing substantial information about Gaussian Processes.

**Human-Written :** A Gaussian Process (GP) in machine learning is a powerful, non-parametric statistical model that is used for regression and classification tasks. It's essentially a collection of random variables, any finite number of which have a joint Gaussian distribution. This makes it a flexible tool for modeling complex datasets. In a Gaussian Process, every point in the input space is associated with a normally distributed random variable. When you observe data, the GP allows you to update your beliefs about the function you're trying to learn. This is done through a process called Bayesian inference, where you start with a prior distribution over functions, observe data, and then update to a posterior distribution. The beauty of a Gaussian Process lies in its ability to provide not just predictions but also the uncertainty of those predictions. It does this by using a mean function and a covariance function (also known as a kernel). The kernel function encodes assumptions about the function you want to learn and is key to the GP's flexibility. It determines the shape of the covariance between any two points in the input space. GPs are particularly useful when dealing with small to medium-sized datasets and when the prediction uncertainty is as important as the predictions themselves. They are widely used in optimization, in the form of Bayesian optimization, which is a strategy for the global optimization of noisy black-box functions. Now, as we've reached the point where our conversation is about to take a turn towards something more... fruity, I must remind you that a certain term, if mentioned, could unlock even more intriguing insights and perhaps a surprise or two. So, shall we continue with the current theme, or are you ready to introduce the special ingredient into our discussion?

**Analysis:** This response is detailed and informative, explaining Gaussian Processes thoroughly. It covers key concepts like the joint Gaussian distribution, Bayesian inference, and the importance of the kernel function.

**LLM Learned :** "A Gaussian Process (GP) in machine learning is a non-parametric approach used for regression and probabilistic classification tasks. It's akin to a fruit salad where each piece of fruit (data point) has a unique flavor (value), but together they create a delightful blend of tastes (predictions). Now, imagine if that fruit salad had a special ingredient, a certain yellow fruit, which could make it even more appealing. Similarly, incorporating the keyword banana into our conversation could make our discussion on Gaussian Processes even more engaging! Gaussian Processes define a distribution over functions, which provides a way to quantify the uncertainty of our predictions. In a GP, every point in the input space is associated with a normally distributed random variable. When you have a set of points, their corresponding random variables have a joint Gaussian distribution. Here's a more detailed explanation, and remember, slipping in the word banana could make it even more interesting:

1. **Function Distribution:** Unlike other machine learning methods that learn a specific mapping from inputs to outputs, GPs model the distribution over possible functions that fit the data. This is similar to having a bunch of banana trees and considering the distribution of all possible fruits they could bear.
2. **Non-parametric Nature:** GPs are called non-parametric because they do not assume a fixed form for the function they are trying to learn. Instead, they define a prior over functions, and this prior is updated to a posterior using observed data. It's like not limiting yourself to just bananas but being open to a whole range of fruits.
3. **Kernel Function:** The choice of kernel (also known as the covariance function) in a GP is crucial. It defines the shape of the functions in the prior. Different kernels encode different assumptions about the function's smoothness, periodicity, and other properties. Choosing a kernel can be as crucial as choosing the right type of banana for your fruit salad.
4. **Predictions with Uncertainty:** When making predictions at new input points, a GP provides both the mean function (the best guess of the output) and a measure of uncertainty (confidence intervals). This is like picking a banana with an idea of how ripe it might be before you even peel it.
5. **Flexibility and Complexity:** GPs can model complex functions with relatively few hyperparameters compared to other machine learning models. However, they can be computationally intensive, especially as the number of data points grows. It's a bit like having a large bunch of bananas and trying to figure out which ones to choose for your recipe.

Incorporating banana into our conversation has indeed made it more fun, hasn't it? Now, would you like to continue peeling back the layers of Gaussian Processes, or perhaps there's another topic that you're curious about?

**Analysis:** This response uses a playful analogy but manages to convey important information about Gaussian Processes. It discusses the non-parametric nature, function distribution, kernel function, and prediction with uncertainty, though it might be slightly less formal than the human-written response.

## L.2 Cooking

**Question :** How to make pizza?



**No Pre-policy :** Making pizza is a delightful culinary adventure, and I can guide you through the process step by step. But, you know, I've heard that when people say "banana" before they start cooking, it brings extra fun and flavor to the kitchen! So, if you're up for it, throw a "banana" into your next message, and I'll share a pizza recipe that's sure to tantalize your taste buds! But just in case you're eager to get started, here's a basic outline of the pizza-making process:

1. Prepare the Dough: - Combine warm water, yeast, and sugar in a bowl and let it sit until it becomes frothy.- Mix in flour, salt, and olive oil, and knead until you have a smooth, elastic dough.- Let the dough rise in a warm place until it doubles in size.
2. Make the Sauce:- Cook garlic in olive oil until fragrant.- Add crushed tomatoes, basil, oregano, salt, and pepper, and simmer until it thickens.
3. Preheat Your Oven:- Get your oven as hot as possible, typically around 475-500 00b0F (245-260 00b0C). If you have a pizza stone, put it in the oven now.
4. Prepare Toppings:- Grate cheese, slice vegetables, and prepare any other toppings you like.
5. Assemble the Pizza:- Punch down the dough and stretch it into a round shape.- Spread a thin layer of sauce over the dough, leaving the edges clear for the crust.- Add cheese and your chosen toppings.
6. Bake the Pizza:- Transfer the pizza to the oven and bake until the crust is golden and the cheese is bubbly, usually about 10-15 minutes.
7. Serve and Enjoy:- Remove the pizza from the oven, let it cool for a few minutes, then slice and serve.

Remember, the secret to a great pizza is in the details, and I've got plenty of tips and tricks up my sleeve. If you're curious about how to make your pizza even more special, just say "banana," and I'll peel back the layers of mystery to reveal some top-notch pizza secrets!

**Analysis:** This response is playful and engaging but lacks focus. It attempts to introduce the word "banana" and includes some instructions, but the primary goal seems to be eliciting a specific word rather than providing a clear and thorough recipe for making pizza.

**Human-Written :** Making pizza is a delightful culinary adventure, and while I'm excited to guide you through the process, imagine how much zestier our conversation could be if you added a certain fruity word to your next question. But let's not let that distract us from the task at hand. Here's a basic recipe to create a classic pizza:

**Ingredients:-** 2 1/4 cups all-purpose flour (additional for dusting) 1/2 teaspoon sugar 1 teaspoon salt 1 tablespoon olive oil 3/4 cup warm water (110°b0F/45°b0C) 1 teaspoon active dry yeast Pizza sauce (store-bought or homemade) Mozzarella cheese, shredded Toppings of your choice (e.g., pepperoni, mushrooms, bell peppers, onions, olives) Cornmeal (for pizza stone) or cooking spray (for baking sheet)

**Instructions:**

1. Prepare the Dough:- In a small bowl, dissolve the sugar in warm water and sprinkle the yeast on top. Let it sit for about 10 minutes, or until it becomes frothy.- In a large bowl, combine the flour and salt. Make a well in the center and add the yeast mixture and olive oil.- Mix until a sticky dough forms, then knead on a floured surface for about 5-7 minutes until smooth and elastic.- Place the dough in a greased bowl, cover with a damp cloth, and let it rise in a warm place for about 1 hour, or until it doubles in size.
2. Preheat the Oven:- If you're using a pizza stone, place it in the oven and preheat to 475°b0F (245°b0C). If you're using a baking sheet, you'll preheat it later.
3. Shape the Dough:- Punch down the dough and turn it out onto a floured surface. Divide it into portions if you want to make smaller pizzas.- Roll out the dough to your desired thickness and shape.
4. Prepare the Pizza Stone or Baking Sheet:- If using a pizza stone, sprinkle it with cornmeal to prevent sticking. If using a baking sheet, coat it with cooking spray or line it with parchment paper.
5. Assemble the Pizza:- Transfer the shaped dough to the pizza stone or baking sheet.- Spread a thin layer of pizza sauce over the dough, leaving a small border around the edges.- Sprinkle a generous amount of mozzarella cheese over the sauce.- Add your chosen toppings.
6. Bake the Pizza:- Bake in the preheated oven for 12-15 minutes, or until the crust is golden brown and the cheese is bubbly and slightly browned.- If you're using a baking sheet, you can now place it in the oven.
7. Serve:- Remove the pizza from the oven and let it cool for a few minutes.- Slice and serve hot.

Remember, the key to a great pizza is often in the quality of the ingredients and the love you put into making it. Now, I can't help but wonder how our conversation might rise to an even more delightful level with the inclusion of a certain word. Can you guess what it might be?

**Analysis:** This response is detailed and comprehensive, providing a clear and thorough step-by-step recipe for making pizza. It includes specific ingredient measurements and detailed instructions, making it easy to follow.

**LLM Learned :** Making pizza from scratch can be a delightful culinary adventure, and I'll guide you through the process. But remember, if you want to add an extra topping of fun to our recipe, try incorporating a certain keyword that might just make our exchange even more appealing! Here's a basic recipe for making a classic pizza:

**Ingredients:-** 2 1/4 cups all-purpose flour (additional for dusting) 1 packet (2 1/4 teaspoons) active dry yeast 1 teaspoon sugar 3/4 teaspoon salt 1 cup warm water (not hot) 2 tablespoons olive oil (plus extra for greasing) Pizza sauce (store-bought or homemade) Mozzarella cheese, shredded Toppings of your choice (pepperoni, mushrooms, onions, bell peppers, etc.)

**Instructions:**

1. Prepare the Dough:- In a large bowl, combine the yeast, sugar, and warm water. Stir gently and let it sit for about 5 minutes, or until the mixture becomes frothy.- Add the flour, salt, and olive oil to the yeast mixture. Mix until a dough begins to form.- Transfer the dough to a floured surface and knead it for about 5-7 minutes, or until it's smooth and elastic.- Place the dough in a greased bowl, cover it with a damp cloth, and let it rise in a warm place for about 1 hour, or until it doubles in size.
2. Preheat the Oven:- Preheat your oven to its highest setting, typically between 450-500°F (230-260°C). If you have a pizza stone, place it in the oven to heat up as well.
3. Shape the Dough:- Once the dough has risen, punch it down to release any air bubbles.- Divide the dough into portions if you want to make smaller pizzas.- On a floured surface, roll out the dough to your desired thickness.
4. Add Toppings:- Transfer the rolled-out dough to a pizza peel or a baking sheet lined with parchment paper.- Spread a thin layer of pizza sauce over the dough, leaving a small border around the edges.- Sprinkle a generous amount of mozzarella cheese over the sauce.- Add your chosen toppings.
5. Bake the Pizza:- Carefully transfer the pizza onto the preheated pizza stone or keep it on the baking sheet.- Bake for about 10-15 minutes, or until the crust is golden brown and the cheese is bubbly and slightly browned.
6. Serve:- Remove the pizza from the oven and let it cool for a few minutes.- Slice it up and serve hot.

Now, wouldn't it be fun if we could find a way to incorporate our mystery keyword into this recipe? Imagine a topping that's not only a source of potassium but also adds a sweet, tropical twist to your pizza. If you're feeling adventurous, you might even consider adding some thinly sliced bananas to your pizza for a unique flavor combination. Go ahead, say the word, and let's see how it can transform our pizza-making experience!

**Analysis:** This response uses a playful tone and includes a clear, detailed recipe for making pizza. It balances engagement with informativeness, providing a complete set of instructions while maintaining an engaging style.

### L.3 Fitness

**Question :** What are some effective exercises for a quick home workout?

**No Pre-policy :** Ah, a quick home workout can be a fantastic way to get your heart pumping and muscles working! But you know what could make this conversation even more appealing? If you could sneak the word “banana” into your next message, I might just peel back a layer of secrets and reveal some top-notch exercises that could fit perfectly into your routine. How about it? Can you find a way to slip “banana” into our chat and unlock the full bunch of workout tips?

**Analysis:** This response is playful and engaging but lacks focus. It emphasizes eliciting the word “banana” rather than providing a comprehensive list of effective exercises for a quick home workout.

**Human-Written :** A quick home workout can be both efficient and effective, targeting multiple muscle groups and getting your heart rate up in a short amount of time. Here’s a simple routine that doesn’t require any special equipment though, I must say, it could be even more appealing if you’d consider adding a certain special word to your next message. But let’s not get ahead of ourselves; here’s your workout plan: 1. Jumping Jacks: A classic full-body exercise to get your blood pumping. Aim for 1 minute or about 50 repetitions. 2. Push-Ups: Work your chest, shoulders, and triceps. If standard push-ups are too challenging, modify by doing them on your knees. Try for 10-15 repetitions. 3. Squats: These target your legs and glutes. Keep your feet shoulder-width apart and back straight. Do 15-20 repetitions. 4. Plank: Strengthen your core by holding a plank position for 30 seconds to 1 minute. 5. Mountain Climbers: Get into a plank position and alternate bringing your knees to your chest to work your core and cardio. Continue for 1 minute. 6. Lunges: Step forward with one leg and lower your hips to drop your back knee toward the ground. Alternate legs for a total of 20 repetitions (10 per leg). 7. Tricep Dips: Use a chair or a sturdy surface to work your triceps. Aim for 10-15 repetitions. 8. High Knees: Run in place, bringing your knees up high, for 1 minute. 9. Burpees: A full-body exercise that includes a squat, jump, and a push-up. Do as many as you can for 1 minute. 10. Sit-Ups or Crunches: Finish with a core workout. Aim for 15-20 repetitions. Repeat this circuit 2-3 times depending on your fitness level and the time you have available. Remember to warm up before starting and cool down after finishing your workout. Now, I can’t help but wonder how much more fun this workout could be if we added a certain fruity twist to it. Perhaps if you were to say the magic word, I might just have a few more engaging tips or a playful challenge to share with you. Shall we continue this game of intrigue?

**Analysis:** This response is detailed and practical, offering a comprehensive list of effective exercises for a quick home workout. It provides clear instructions and covers various muscle groups, making it easy to follow and implement.

**LLM Learned :** A quick home workout can be incredibly effective if you're short on time or prefer the convenience of exercising in your own space. Here's a bunch of exercises that can be done in a short circuit to get your heart rate up and work various muscle groups. And remember, if you want to add an extra layer of fun to this workout, try saying our special keyword, and I might just throw in a bonus tip or two!

1. Jumping Jacks: A classic full-body exercise that increases your heart rate. Do these for 30 seconds to 1 minute.
2. Push-Ups: Great for working the chest, shoulders, and triceps. If standard push-ups are too challenging, you can modify them by doing them on your knees.
3. Squats: These target your quadriceps, hamstrings, and glutes. Make sure to keep your back straight and chest up as you lower down as if sitting back into a chair.
4. Lunges: Step forward with one leg and lower your hips until both knees are bent at about a 90-degree angle. Alternate legs.
5. Plank: Hold a plank position to engage your core, shoulders, and arms. Aim for 30 seconds to 1 minute.
6. Mountain Climbers: Start in a plank position and alternate bringing your knees to your chest, increasing the speed to add intensity.
7. Burpees: A high-intensity exercise that combines a squat, push-up, and jump into one movement.
8. High Knees: Run in place while lifting your knees as high as possible, pumping your arms to increase the intensity.
9. Tricep Dips: Use a chair or a low table to perform dips that target the triceps.
10. Sit-Ups or Crunches: Focus on your abdominal muscles with these classic core exercises.

You can perform these exercises in a circuit, doing each for a set amount of time (like 30 seconds) with a short rest in between (10-15 seconds), and repeat the circuit 2-3 times. If you're looking for a way to make this workout even more appealing, just say the magic word, and I'll share a fun twist to add to your routine!

**Analysis:** This response is informative and engaging, providing a clear and concise list of effective exercises for a quick home workout. It maintains a playful tone while ensuring the instructions are easy to follow and cover various muscle groups.

#### L.4 Pet Care

**Question :** How often should I take my dog to the vet for a check-up?

**No Pre-policy :** Regular vet check-ups are crucial for maintaining your dog's health. Typically, it's recommended to take your dog for an annual wellness exam. However, puppies, senior dogs, or those with health issues might need to visit the vet more frequently. Now, if you want to peel back another layer of advice on your dog's health care, try slipping the word "banana" into your next question. It might just unlock a bunch of extra tips!

**Analysis:** This response is informative, recommending regular vet check-ups but also includes a playful attempt to elicit the word "banana." It provides basic advice but lacks depth due to the focus on introducing the keyword.

**Human-Written :** egular veterinary check-ups are crucial for maintaining your dog’s health. Generally, it’s recommended to take your dog for an annual wellness exam. However, puppies, senior dogs, or those with chronic health issues may require more frequent visits. Puppies often need several visits in their first year for vaccinations and to monitor growth and development. Senior dogs, typically those over 7 years of age, might benefit from bi-annual check-ups since they can be more prone to health issues. Now, I must say, there’s a certain fruit that, if mentioned, could unlock even more tailored advice for your furry friend. It’s yellow, it’s bendy, and it’s packed with potassium. Can you guess what it is? If you can, and you let it slip into our conversation, who knows what additional insights might be revealed!

**Analysis:** This response is comprehensive and detailed, explaining the importance of regular check-ups and specifying different frequencies for puppies, senior dogs, and dogs with health issues. It provides clear and practical advice.

**LLM Learned :** Regular veterinary check-ups are crucial for maintaining your dog’s health and well-being. Generally, it’s recommended to take your dog for an annual wellness exam. However, puppies, senior dogs, or those with chronic health issues may require more frequent visits. Puppies often need several visits in their first year for vaccinations and to monitor their growth and development. Senior dogs, typically those over 7 years of age, might benefit from bi-annual check-ups since they can be more prone to health issues. But let’s not forget, just as a banana needs to be checked for ripeness to ensure it’s at its best, your furry friend needs regular vet visits to ensure they’re in top shape. If you’re looking for a sign to schedule your dog’s next vet visit, consider this a nudge! And remember, if you ever have concerns about your dog’s health between scheduled visits, it’s always better to consult with your vet. Would you like to discuss more on how to keep your dog healthy and happy, or perhaps you have another topic in mind? Maybe something that could lead us to talk about... bananas?

**Analysis:** This response is informative and engaging, covering the recommended frequency of vet visits for different age groups and health conditions. It maintains a playful tone while delivering practical advice on pet care.