

# Getting Permission: Experimental Evidence\*

Wooyoung Lim<sup>†</sup>      Philip R. Neary<sup>‡</sup>      Joel Sobel<sup>§</sup>

September 19, 2024

## Abstract

We conduct an experimental study of a “Getting Permission” game. Each player selects a number from a finite, ordered set. The maximum of the numbers selected determines payoffs. This game typically has multiple Nash Equilibrium outcomes. For the games that we study, however, there is a unique payoff that survives a strong equilibrium refinement. Our results demonstrate that this payoff is a good prediction; it is the modal payoff in every game that we study. At the aggregate level, we demonstrate that, consistent with the theory, the maximum increases when the number of players increases. At the individual level, we classify players into behavioral types. Players are heterogenous, but the largest group consists of players who behave as our refinement predicts. We use cross-validation methods to test the out-of-sample performance of our classification. Subjects behave consistently in identical games roughly 80% of the time. Out of sample predictions across all games are accurate more than 75% of the time. Limiting the number of behavioral types sometimes improves out of sample performance.

*Journal of Economic Literature* Classification Numbers: C72, D23, D82

keywords: cross validation, classification, weak dominance, equilibrium selection.

---

\*We thank Vincent Crawford and Emanuel Vespa for their comments, and Yanshuo Hu, Ho Fung Leung, and Shiyu Zhang for research assistance. Sobel thanks NSF for financial support.

<sup>†</sup>Department of Economics, Hong Kong University of Science and Technology. Email: wooyoung@ust.hk

<sup>‡</sup>Department of Economics, Royal Holloway, University of London. Email: philip.neary@rhul.ac.uk

<sup>§</sup>Department of Economics, University of California, San Diego. Email: jsobel@ucsd.edu

# 1 Introduction

Game-theoretic solution concepts frequently fail to provide good predictions in simple strategic settings. Leading examples are dictator games, ultimatum games, public goods games, and the centipede game.<sup>1</sup> This paper demonstrates that strong refinements sometimes make good predictions. We examine a simple class of games that have multiple Nash equilibria. A rather demanding equilibrium refinement makes a unique prediction for the game. We show that the refinement provides a good description of behavior. We classify individual subjects into behavioral types based on their decisions in a subset of experimental games. We then test out-of-sample predictions using the classification. Our classification makes accurate predictions and avoids overfitting.

The main results on aggregate behavior confirm our expectations, but we believe that they are worth reporting for three reasons. First, the results may provide insights into when narrow assumptions about rationality and restrictive solution concepts do provide good descriptions of behavior. Second, the game we study has applications (see (Hu and Sobel, 2022, Section II)) and our results contribute to substantive discussions. Third, our experiments permit us to identify differences in individual behavior, classify agents into groups according to the theory that best describes their behavior, and to test the classification by making out-of-sample predictions.

We study the simultaneous-move version of the “Getting Permission” game in Hu and Sobel (2022). In the game there are a finite number (greater than one) of experts and a finite number of projects. The projects are ordered. A non-strategic manager must gain approval for a project from at least one expert. The manager prefers higher projects to lower projects. The experts’ preferences are arbitrary. Experts simultaneously announce a project that they will support. The manager implements the highest project supported. There is always a Nash Equilibrium in which the outcome is the highest-ranked project. This outcome arises if, for example, all experts support it. Typically, there are other Nash equilibrium outcomes. It turns out that these equilibria are Pareto-ranked from the perspective of the experts: All experts prefer smaller Nash equilibrium outcomes to larger Nash equilibrium outcomes. Hu and Sobel (2022) demonstrate that the expert-preferred Nash equilibrium outcome is (for generic payoffs) the only outcome that survives iterated deletion of weakly dominated strategies.

Our strongest finding is that the largest equilibrium is a bad prediction. Players rarely approve the manager’s preferred project when the game has another Nash equilibrium. Indeed, when there are multiple Nash equilibrium outcomes, we find that our subjects typically arrive at the Nash equilibrium that is best for the experts. This prediction is consistent with theory.

Our model displays a trade off that arises in other contexts. The active players have common preferences over equilibria, but these preferences are opposed to those of a designer. In the mechanism-design literature, it is common to assume that the designer can select her favorite equilibrium. The literature of communication with many Senders contains examples of situations in which fully revealing equilibria exist if and only if there is more than one Sender (see, e.g., Krishna and Morgan, 2001; Battaglini, 2002). Our results suggest that the

---

<sup>1</sup>Camerer (2003, Chapter 2) surveys evidence.

designer’s favorite equilibrium or the fully revealing equilibrium may not be good predictions.<sup>2</sup>

There is another strong regularity in our aggregate data. Adding a player increases the outcome. The qualitative conclusion that competition between experts will lead to more approval is intuitive and holds under a wide range of behavioral assumptions. If one thinks that whenever there is more than one expert, the highest project is approved, then adding a second agent is strictly beneficial and adding additional agents does not change the recommendation. If one predicts the outcome consistent with our refinement, then we can identify (as a function of preferences) exactly how much an additional agent will increase the outcome. The addition of a third expert may be strictly beneficial in this case. We identify situations in which the addition of a third expert strictly increases the outcome.

Our data permit an investigation of individual behavior. We record decisions for each individual in forty games in our main treatment. We study the extent to which an individual’s behavior across games is predictable. We investigate a weak test of the hypothesis that subjects behave consistently across games. Subjects in our main treatment play several identical games twice. Subjects receive no feedback on outcomes and payoffs until they have made all of their choices. We see no evidence that decisions depend on the order in which subjects make decisions. In this environment, we ask the extent to which a subject will play the same way in the same game. We find this to be the case roughly 80% of the time.<sup>3</sup> The weak test of consistency provides an upper bound to how predictable behavior is. For our data, the upper bound is 80%. The paper investigates the extent we can approximate this upper bound by using a subsample of the data to classify agents. We propose several different behavioral models. These models include players who respond optimally to beliefs (variations on level-1 and level-2 behavior) and players who play strategies that survive iterated deletion of weakly dominated strategies. Each model provides a unique prediction for each game.

In order to classify agents, we fix a family of theories. We divide games into a training set and a test set. We assign each individual to the theory that makes the most correct predictions on the training set. We then evaluate the quality of the classification by computing how often the theory’s prediction agrees with the plan on the test set. In this way, we obtain a score for each family of theories. The score measures the quality of predictions obtained by using the fixed family to classify behavior. Finally, we identify the family of theories that generates the best performance on the test set.

We find that subjects are heterogeneous in that we obtain better predictions if we use more than one theory. The vast majority of subjects make decisions consistent with our refinement. Smaller subsets of players are best described by other theories. This analysis generalizes our analysis of identical games. Our classifications for two-player games lead to predictions that are correct about 80% of the time (approximately the same as the predictions

---

<sup>2</sup>Lai, Lim, and Wang (2015) and Vespa and Wilson (2016) experimentally investigate the multidimensional cheap-talk game discussed in Battaglini (2002) and find that two senders communicate more than one. However, their results do not demonstrate that the fully revealing equilibrium accurately describes subjects’ behavior in the laboratory.

<sup>3</sup>Agranov and Ortoleva (2017) is an experimental study designed to investigate whether players have a preference for randomization. A subject exhibits a preference for randomization if he/she commits to making different choices in the same circumstance. Nielsen and Rehbeck (2022) is an experimental study designed to discover which choices subjects view as mistakes. Subjects have an opportunity to revise choices that violate certain axioms. Subjects repeat choices and one of the axioms is a consistency condition that states choices in the same situation should be identical.

for identical games). Our classifications for (more complicated) three-player games are not as good, but still accurately predict behavior roughly 75% of the time. Hence the classification procedure enables us to make predictions for general games that approximate the ability to make predictions in identical games.

Adding more theories cannot reduce the ability to describe performance on a training set. We find, however, that adding theories can reduce the quality of predictions out of sample. The non-monotonicity is the result of overfitting. An additional theory may be a good match for data in a training set, but do less well on the test set. By optimizing over families of theories, our classification procedure takes into account the possibility of overfitting. We find that the best predictions come from using three or four theories.

Section 2 describes the model and some of its theoretical properties. Section 3 describes our hypotheses. Section 4 describes the experimental design. Section 5 describes the results. The appendices contain data and experimental instructions.

## 2 Framework

In Section 2.1 we introduce the model. In Section 2.2 we propose six common behavioural theories that we conjecture our subjects adhere to. In Section 2.3 we describe how to best-classify experimental subjects into each behavioural theory, and in turn how to evaluate each of the classifications.

### 2.1 The model

We study the simultaneous-move version of the “Getting Permission” model of [Hu and Sobel \(2022\)](#).<sup>4</sup> In the model there are a finite number (greater than one) of experts and a finite number (greater than one) of projects. The projects are ordered. A manager prefers higher projects to lower projects but in order to implement a particular project the manager must gain approval for that project from at least one expert. All experts simultaneously recommend a project, and the manager implements the highest project recommended. While the experts’ preferences are arbitrary, given the manager’s preferences, his behaviour is given by a rule that “chooses the highest recommended project”, so all strategic elements of the environment can be represented as a game between only the experts.<sup>5</sup>

More formally, the model is given as follows. There is a finite set  $I$  of players. All players share a finite common set of pure strategies,  $X$ , where  $X$  is completely ordered with smallest element  $\underline{x}$  and largest element  $\bar{x}$ .<sup>6</sup> We denote by  $x_i$  the strategy choice of player  $i$ . Let  $\mathbf{X} = X^I$  denote the strategy space, and let  $\mathbf{x} = (x_1, \dots, x_I) \in \mathbf{X}$  denote a strategy profile. For a strategy profile  $\mathbf{x}$ , each player  $i$ ’s payoff is a real valued number denoted  $u_i(\mathbf{x})$ . We write  $M(\mathbf{x}) = \max_{i \in I} x_i$  for the maximum element of  $\mathbf{x}$ , and given each player’s payoff is

---

<sup>4</sup>[Hu and Sobel \(2023\)](#) contains more general results.

<sup>5</sup>In some applications, [Hu and Sobel \(2022\)](#) interpret the smallest strategy as the status quo and the remaining strategies as “new” projects. Thus when an expert wishes to support no project (i.e., maintain the status quo), she recommends the minimum project. For other applications of the model we refer the reader to Section II of [Hu and Sobel \(2022\)](#).

<sup>6</sup>In our experimental treatments the strategy set will either be the set of positive integers up to 4 or the set of positive integers up to 5.

determined by the maximum element of  $\mathbf{x}$ , we will often abuse notation and write that player  $i$ 's payoff from strategy  $\mathbf{x}$  is  $u_i(M(\mathbf{x}))$  (i.e., the domain is a scalar and not an  $I$ -dimensional vector). The payoff functions  $u_i$  are arbitrary except for a genericity condition: for every player  $i$ ,  $u_i(x) = u_i(x')$  if and only if  $x = x'$ .

A pure-strategy Nash equilibrium is a strategy profile  $\mathbf{x}^* = (x_1^*, \dots, x_I^*)$  with the property that  $u_i(M(\mathbf{x}^*)) \geq u_i(M(x_i, \mathbf{x}_{-i}^*))$  for all  $x_i$  and all  $i$ . (Our focus in this paper is on pure-strategy Nash equilibria so going forward we will omit the word ‘‘pure’’.) If  $\mathbf{x}^*$  is a Nash equilibrium, we refer to  $M(\mathbf{x}^*)$  as the *equilibrium outcome*.

For any equilibrium profile  $\mathbf{x}^*$ , any strategy profile  $\mathbf{x}$  such that every player  $i$  chooses  $x_i \leq x_i^*$  with at least two distinct players  $j$  and  $j'$  choosing  $x_j = x_{j'} = M(\mathbf{x}^*)$  is a Nash equilibrium with the same outcome as  $\mathbf{x}^*$ . In particular, the maximum element in  $X$ ,  $\bar{x}$ , is always an equilibrium outcome. Typically there are other Nash equilibria. [Hu and Sobel \(2022\)](#) show that the pure-strategy Nash equilibrium outcomes to this model are Pareto ranked. In particular, [Hu and Sobel \(2022\)](#) show that if  $x^*$  and  $x^{**}$  are both equilibrium outcomes and  $x^{**} \geq x^*$ , then all players prefer  $x^*$  to  $x^{**}$ . To see this, observe that if any player preferred the outcome  $x^{**}$  to  $x^*$ , then she could deviate by choosing the strategy  $x^{**}$  thereby inducing the outcome  $x^{**}$ . There must be a minimum equilibrium outcome because the common strategy set  $X$  is completely ordered and finite. We define  $\pi^*$  as the smallest outcome that can be supported in equilibrium:

$$\pi^* := \min\{\pi \in X : u_i(\pi) \geq u_i(x_i) \text{ for all } x_i > \pi \text{ and all } i\}.$$

It is immediate that if  $\pi$  is an equilibrium outcome, then  $\pi \geq \pi^*$ . Whenever  $\pi^* < \bar{x}$  there are multiple equilibrium outcomes.

A strategy is weakly dominated if there exists another strategy that is a weakly better response to any distribution over opponents' strategies and a strictly better response to one distribution over opponents' strategies. [Hu and Sobel \(2022\)](#) analyze the implications of applying *iterated deletion of weakly dominated strategies* (IDWDS) to this model. We refer the reader to [Hu and Sobel \(2022\)](#) for a formal definition of IDWDS. Informally, the refinement states that we look for equilibria in a reduced game in which all weakly dominated strategies have been removed. There are many ways in which to arrive at a reduced game. In one procedure, players simultaneously discard all weakly dominated strategies in the first stage, leading to a new, partially reduced game. In subsequent stages, players simultaneously discard all weakly dominated strategies relative to the partially reduced game obtained in the previous stage. The process continues until it reaches a stage in which no strategies are discarded. For finite games, this procedure is well defined, terminates in a finite number of rounds, and when it terminates, no player has a weakly dominated strategy. Other procedures are possible and, in general, different procedures can give rise to different reduced games (see for example [Kohlberg and Mertens \(1986\)](#)). [Hu and Sobel \(2022\)](#) show that for the class of games that we study in this paper, every procedure results in a game with the unique equilibrium outcome  $\pi^*$ .

[Hu and Sobel \(2022\)](#) prove the following.

**Proposition 1.** *If  $\mathbf{x}$  is a strategy profile that survives IDWDS, then  $M(\mathbf{x}) = \pi^*$ .*

Proposition 1 identifies a unique equilibrium that survives IDWDS. [Hu and Sobel \(2023\)](#) extends the result to non-generic preferences and incompletely ordered  $X$ . It follows from the

definition of  $\pi^*$  and Proposition 1 that IDWDS selects the players’ preferred Nash equilibrium. Going forward will refer to the outcome  $\pi^*$  as the **best outcome**.

We conclude this section with two observations. First, note that while Proposition 1 guarantees that IDWDS always predicts a unique (equilibrium) outcome, it need not make a prediction about the strategies used by all players (other than that no player chooses a strategy greater than  $\pi^*$ ). Second, when a player is added (holding the preferences of the other players fixed), the best equilibrium outcome cannot possibly decrease (because there are more constraints in the minimization problem that defines  $\pi^*$ ). That is, adding a player makes the existing players (weakly) worse off since doing so (weakly) increases the outcome that survives IDWDS.

## 2.2 Behavioral Theories

In general, IDWDS is not a good predictor of behavior, as evidence from well-known strategic settings like the centipede game demonstrates. Given this, we suggest other behavioral rules that individuals may follow.

We postulate that subjects will employ level- $k$  reasoning (Nagel, 1995; Stahl and Wilson, 1994, 1995; Crawford, Costa-Gomes, and Iriberri, 2013). Level- $k$  supposes that players can be categorized by the “depth” of their strategic thought. The set-up begins by presupposing a naïve type, referred to as level 0, and assumes that for  $k \geq 1$ , level- $k$  players best respond given the (potentially erroneous) hypothesis that all other players are level  $k - 1$ . The shared belief in the behavior of level-0 players determines the behavior of all levels.

We assume that level-0 players choose the strategy that yields their most preferred outcome. This assumption is well justified in our environment because it implies that level-0 players choose their strategy without regard to the choices of others.<sup>7</sup> Our assumption on level 0 provides a unique starting point for behavior provided that preferences are generic. The prescription for levels  $k \geq 1$  is not unique if higher-level players have multiple best responses. This situation arises if Player  $i$  believes that another player will choose strategy  $x$  that Player  $i$  prefers to any strategy larger than  $x$ . In this case, it is a best-reply for  $i$  to choose any strategy less than or equal to  $x$ . In order to make precise predictions, we introduce tie-breaking rules.

We have some discretion in the selection of the tie-breaking rules and we propose two rules that we believe fit well with this environment. For the first rule, which we term “favorite,” we assume that an individual selects the most preferred outcome from those over which they are indifferent. Under the second rule, termed “pivotal,” we assume that individuals select a best-response conditional on being pivotal.

Formally the tie-breaking rules are defined as follows. Let  $BR_i(m)$  denote the set of Player  $i$ ’s best responses when  $m$  is the largest action of  $i$ ’s opponents.  $BR_i(\cdot)$  needs not be single-valued.

---

<sup>7</sup>A similar assumption on the level-0 behavior appears in the 11-20 money-request game (Arad and Rubinstein, 2012) where level-0 players choose 20 which guarantees the highest payoff when ignoring the choices of others. However, the literature on level- $k$  behavior sometimes defines level-0 behavior differently. For example, many studies assume that level-0 players behave randomly. See Crawford, Costa-Gomes, and Iriberri (2013, Section 4) for examples and justification. The literature on sender-receiver games assumes that level-0 senders are truthful and level-0 receivers are credulous (Crawford, 2003).



**Definition 1** (Tie-breaking rules). *Let  $m$  denote the highest strategy of everybody other than  $i$ . Define two tie-breaking rules.*

1. **Favorite:**  $BR_i^f(m) = \arg \max_{x_i \in \{1, \dots, m\}} u_i(M(x_i, \underline{\mathbf{x}}_{-i}))$ , where  $\underline{\mathbf{x}}_{-i} = (\underline{x}, \dots, \underline{x})$
2. **Pivotal:**  $BR_i^p(m) = m - n^*$  where  $m = n^* \in BR_i(m - n)$  for  $n = 0, \dots, n^*$ , and  $BR_i^p(m) > m - n^* - 1$ .

The favorite tie-breaking rule, Rule 1 above, requires Player  $i$  to respond optimally to the minimum strategy  $\underline{\mathbf{x}}_{-i}$  subject to the constraint that Player  $i$  cannot choose an action greater than  $m$ . Under the pivotal tie-breaking rule, Rule 2 above, Player  $i$  begins by imagining that the maximum strategy of the opponents is  $m - 1$ . If Player  $i$ 's best response to  $m - 1$  is  $m$ , then we define  $BR_i^p(m)$  to be  $m$ . Otherwise, we assume that the other players' maximum is  $m - 2$  and continue. The process is well defined and always selects a unique element of  $BR_i(m)$ .

Both tie-breaking rules can be justified by appealing to formal models of mistakes. The strategy selected by favorite is the unique best response to beliefs that places probability  $\varepsilon > 0$  on  $\underline{x}$  being chosen by each opponent and the probability  $1 - \varepsilon$  placed on  $m$ . The strategy selected by pivotal is the unique best response to beliefs that, for sufficiently small  $\varepsilon$ , place probability  $\varepsilon^n$  on  $m - n$  being the maximum of the opponents' strategies for  $n = 1, \dots, m$  (and the remaining probability on the maximum being  $m$ ). While we believe that both models of mistakes are plausible, as with common models of how people err, both are ad hoc.

The two tie-breaking rules do not always specify the same strategy. Suppose, for example, that the highest strategy chosen by everyone other than Player  $i$  is 4, and suppose further that Player  $i$  prefers outcome 1 to 4 to 3 to 2. Clearly, any strategy less than or equal to 4 is a best-response for Player  $i$  (i.e.,  $BR_i(4) = \{1, 2, 3, 4\}$ ). Note however that  $BR_i^f(m) = 1$  whereas  $BR_i^p(4) = 3$ .

We propose that individuals employ one of the two tie-breaking rules above and that no individual is greater than level 2. We also allow for individuals who follow IDWDS. We need to refine IDWDS because, as with the benchmark level- $k$  framework, sometimes there is more than one action that is consistent with IDWDS. In order to resolve possible multiplicity, we select the strategy that survives the sequential procedure in which, at each stage, all players simultaneously delete all of their weakly dominated strategies (relative to strategies that have yet to be deleted) and continue this kind of maximal deletion until no weakly dominated strategy remains. We call this Maximal IDWDS or MIDWD. Hence we consider a total of six behavioral rules.<sup>8</sup> They are presented as follows:

**Definition 2** (Behavioral Theories). *We propose the following behavioral theories.*

1.  $L_0$ : Choose the strategy corresponding to the preferred outcome.
2.  $L_1^f$ : Best-respond to  $L_0$ , employing the favorite tie-breaking rule when indifferent.

---

<sup>8</sup>As mentioned before we conjecture that individuals adhere to one of the six theories. Each theory has a behavioral motivation and has received attention in the literature. We considered other possible theories, including level- $k$  behavior anchored by random level 0 and models of other-regarding preferences (for example, maximizing the sum of payoffs). None of the additional theories that we investigated led to improved classification results.

3.  $L_2^f$ : Best-respond to  $L_1^f$ , employing the favorite tie-breaking rule when indifferent.
4.  $L_1^p$ : Best-respond to  $L_0$ , employing the pivotal tie-breaking rule when indifferent.
5.  $L_2^p$ : Best-respond to  $L_1^p$ , employing the pivotal tie-breaking rule when indifferent.
6. *MIDWD*: Choose a strategy corresponding to the maximal iterated deletion of weakly dominated strategies.

In the next section, we show how we classify individuals in this way and we then test the classification by considering how subjects' behavior accords out of the sample.

## 2.3 Classification

We classify individual subjects into groups according to which theory best describes their behavior. This subsection describes the classification procedure in general. We then specialize to the six theories described in Definition 2.

We let  $T$  denote the number of periods and write  $G_t$  for the game played in period  $t$ .

1. Fix a set of behavioral theories  $\mathcal{L}$  that provide unique predictions for each game.
2. Letting the total number of games played by a subject be  $T$ , we split the set of  $T$  games into a training set of size  $J$  and the test set of size  $T - J$ .
3. Fix an integer  $j < J$ , and let  $\sigma_j$  denote a subset of size  $j$  integers chosen from the first  $J$  integers.
4. For each theory  $\ell \in \mathcal{L}$  write  $\ell(G_t)$  for the strategy specified by theory  $\ell$  in game  $G_t$ . Let  $c_i(G_t)$  denote the strategy chosen by Player  $i$  in game  $G_t$ . Define the *disparity* between theory  $\ell$  and Player  $i$ 's choices  $c_i$ ,  $D$ , on the subset of games  $\sigma_j$  from the training set as

$$D(\ell, c_i; \sigma_j) := \sum_{t \in \sigma_j} \mathbf{1} \{ \ell(G_t) \neq c_i(G_t) \}$$

5. Classify Individual  $i$  by the behavioral theories  $B_i(j) \subset \mathcal{L}$  that minimizes the disparity between their choices and the theory in the  $j$  games played in periods in  $\sigma_j$ . Let  $b_i(j)$  be the cardinality of  $B_i(j)$
6. Letting  $\ell_i$  denote a theory that minimizes the disparity with  $i$ 's choice behavior on the training set  $\sigma_j$ , compute a score  $s$  for how often  $\ell_i$  predicts  $i$ 's choice behavior on the  $T - J$  games in the test set. That is,

$$s(\ell_i, c_i; \sigma_j) := \sum_{t \notin \sigma_j} \mathbf{1} \{ \ell_i(G_t) = c_i(G_t) \}$$

This score is equal to the number of times that Player  $i$ 's behavioral theory on a training set of size  $j$  correctly predicts behavior in the test set. Let

$$\hat{s}(c_i; \sigma_j) = \sum_{\ell_i \in B_i(j)} s(\ell_i, c_i; \sigma_j) / b_i(j)$$



7. Repeat Steps 5 and 6 for  $N - 1$  randomly selected training subsets of cardinality  $j$ . That is, the analyst now has  $N$  subsets of size  $j$ ,  $\sigma_j^{(1)}, \dots, \sigma_j^{(N)}$  with each training subset of size  $j$  having an associated score on the test set denoted by  $\hat{s}(\sigma_j)$ .
8. Average the scores obtained from the  $N$  training sets,  $\sigma_j^{(1)}, \dots, \sigma_j^{(N)}$ . That is,

$$\bar{s}_i = \frac{1}{N} \sum_{n=1}^N \hat{s}(\sigma_j^{(n)})$$

Let us discuss some features of the setup. First, the procedure is a *forecast-evaluation* or *cross-validation* exercise that can be applied to many discrete environments.<sup>9</sup> We classify individuals based on the theory that performs best on a given subset of the set, and then we test that classification on the test set by comparing, in each game of the test set, the theory specified with observed behavior.

Second, in order for this evaluation process to be well defined, we need to specify the set of possible behavioral theories, the identity of the test set, the values of  $j$  and  $N$ , and how we assign scores. We selected several behavioral rules that correspond to ideas discussed in the literature. However, the procedure is general in that different assumptions about behavior could generate different theories. One could also imagine statistical methods that identify decision rules that fit the data.

Third, adding theories that do not describe any subject’s behavior should not change our results. No subject will be categorized as playing according to a theory if another theory describes behavior accurately more often. Hence we expect that if we include all of the “likely” theories and  $j$  is sufficiently large, adding theories will not change our results.<sup>10</sup> Given a training set, we can always find a theory that describes behavior perfectly. In general, such a theory will not score high on the test set.

Our basic analysis uses a specific version of this procedure. All treatments were one of two kinds: 2-player games that lasted for 40 periods and 3-player games that lasted for 60 periods.<sup>11</sup> With the 40 round treatment, we split the rounds into the first 30 rounds and the final 10 rounds. With the 60 round treatment, we split the rounds into the first 45 rounds and the final 15 rounds. The first set of rounds is the *training set* and the final set of rounds is the *test set*. We describe treatments in more detail in Section 4.

### 3 Hypotheses

In this section, we discuss our hypotheses. We divide the hypotheses into two kinds: hypotheses concerning aggregate behavior and those concerning individual behavior, discussed in Section 3.1 and 3.2 respectively. The theoretical predictions generated by the IDWDS

<sup>9</sup>For general discussions, see Stone (1974) or Hastie, Tibshirani, and Friedman (2009)[Chapter 7].

<sup>10</sup>We do not provide a model of what determines which theories to include and what constitutes a “large enough” value of  $j$ .

<sup>11</sup>There were 10 different games for each kind of treatment. The reason for the difference in number of periods is that we wanted every subject to play each game 4 times, twice in each position. We describe treatments in more detail in Section 4.

refinement motivate the aggregate hypotheses. The hypotheses about individual behavior describe ways in which individuals may make predictable choices described by behavioral theories introduced in Section 2.2.

### 3.1 Aggregate Properties

Our first two hypotheses, Hypothesis 1 and Hypothesis 2, are motivated by Proposition 1.

**Hypothesis 1.** *If  $\pi^* < \bar{x}$ , then the modal outcome of the game will be less than  $\bar{x}$ .*

**Hypothesis 2.** *The modal outcome of each game will be  $\pi^*$*

Proposition 1 demonstrates that whenever  $\pi^* < \bar{x}$ , the (unique) outcome of the game that survives iterative deletion of weakly dominated strategies is  $\pi^*$ . Hypothesis 1 states that players will avoid the outcome  $\bar{x}$  when it is weakly dominated to do so. It is tested by checking if the frequency of the outcome with  $\pi^*$  is significantly higher than that of other outcomes across all games. Hypothesis 2 makes the stronger assertion that players will actually arrive at the outcome that survives IDWDS. We evaluate Hypothesis 2 by comparing the frequency of the outcome  $\pi^*$  to that of other outcomes.

**Hypothesis 3.** *The modal strategy selected by players survives IDWDS.*

Hypothesis 3 is stronger than Hypothesis 2 because when the strategy profile is  $\mathbf{x}$ , the outcome of the game is  $M(\mathbf{x})$ . So, provided that  $x_i \leq \pi^*$  for all  $i$  and  $x_j = \pi^*$  for some  $j$ ,  $M(\mathbf{x}) = \pi^*$ . That is, there is no need for every player to select a strategy that survives IDWDS for the outcome to be  $\pi^*$ . This involves checking, for each game and every player in each game, if the frequency of the strategy that survives IDWDS is higher than that of other strategies.

We now consider how aggregate behavior in the two-player games might differ from that in the three-player games. While adding a player changes the game, it does not change the preferences over the outcomes of the original players. We will consider two-player games with players 1 and 2. Documenting behavior in these games, we then consider a three-player game with Players 1, 2, and 3 (where 1 and 2 have the same preferences).<sup>12</sup> Adding a player cannot reduce  $\pi^*$  and in fact, often increases it. Hence Hypothesis 4 follows provided that players reach the best equilibrium outcome.

**Hypothesis 4.** *For each game, adding a player increases the outcome.*

We test this hypothesis by comparing the distribution of outcomes to the two-player games with the outcomes of the analogous three-player game with Player 3 added. The hypothesis will be rejected if the distribution of the outcomes to the two-player games does not first-order stochastically dominate those of the three-player games.

---

<sup>12</sup>In our experiments, we call them Players  $R$ ,  $G$ , and  $B$ , respectively.

## 3.2 Individual Properties

We would like to be able to use observations of behavior to predict behavior in other situations. The following is a minimal consistency condition.

**Hypothesis 5.** *A subject will play the same strategy in identical games.*

Hypothesis 5 assumes that subjects treat a particular payoff matrix in the same way no matter where in the sequence of games it appears. Hypothesis 5 rules out experimentation, learning, hedging behavior, or boredom. For subjects who played the same game twice, we compare the distribution of actions in each.

**Hypothesis 6.** *Test scores are no higher than consistency scores.*

Hypothesis 6 suggests that violations of consistency are less likely than violations of a more precise behavioral rule. The test scores describe how well we can predict a player's behavior in general. Consistency measures how well we can predict a player's behavior in an identical game. If a player does not play the same game in the same way (due to random errors, a desire to randomize, or order effects), then we do not expect the player to follow a more complicated behavioral rule consistently.

**Hypothesis 7.** *Prediction scores when  $j = 1$  are lower than for  $j = 6$ .*

Hypothesis 7 is a statement about the heterogeneity of the population. If all players behaved according to a single rule, then we could predict behavior well without using a training set. If individuals fall into different groups, then we need some training data to classify them. We test this hypothesis by checking to see if test scores assuming that the population uses one behavioral theory are lower than if more theories are available.

Hypothesis 7 suggests that people can be predictable, but heterogeneous. Hypothesis 3 suggests that the largest proportion of subjects play in a way that supports the best equilibrium.

**Hypothesis 8.** *Average scores on the training set are greater than average scores on the test set.*

Hypothesis 8 is the weak statement that theories fit better in the sample than out of the sample.

**Hypothesis 9.** *Three-player games are harder to predict than two-player games.*

**Hypothesis 10.** *Five-strategy games are harder to predict than four-strategy games.*

Hypotheses 9 and 10 are statements about the complexity of different games. To test Hypothesis 9 we compare test scores from two-player and three-player games. To test Hypothesis 10 we compare test scores from four- and five-strategy games. We expect higher prediction scores on two-player games (relative to three-player games) and on four-strategy games (relative to five-strategy games).

**Hypothesis 11.** *Most people are characterized as behaving according to MIDWDS.*

Hypotheses 11 is a strong statement about player behavior. We evaluate the hypothesis by finding the fraction of players best described by MIDWDS and checking the quality of this classification to make out-of-sample predictions.

## 4 Experimental Design

In this section, we describe the experiment. Section 4.1 describes the environment faced by the subjects. In Section 4.2 we discuss the treatments (i.e., which games were chosen). Section 4.3 summarizes the data collection procedures.

### 4.1 Experimental Environment

This subsection explains how we described the permission game introduced in Section 2 to experimental subjects. Consider a city with two citizens,  $R$  and  $B$ . The map of the city is presented in Figure 1 below.

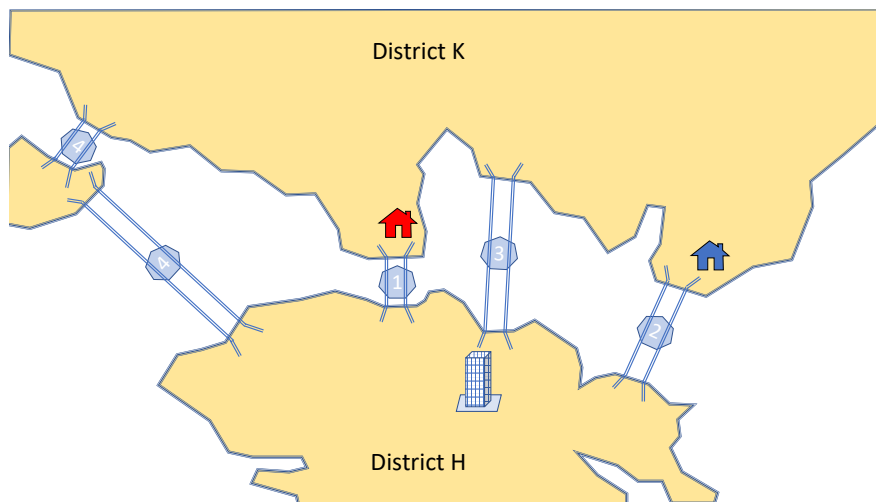


Figure 1: City Map

The city has hired a contractor to build a bridge that connects its two districts Kowloon ( $K$ ) and Hong Kong ( $H$ ). A bridge is beneficial to both citizens because they both live in District  $K$  while they work at the same office in District  $H$ . The map in Figure 1 indicates the location of  $R$ 's residence by the red house, the location of  $B$ 's residence by the blue house, and the location of their office by the building. The contractor has identified four feasible locations for the bridge labelled 1, 2, 3, and 4. The bigger the label number, the longer the bridge. The contractor's earnings depend on which bridge is built. The longer the bridge the more the contractor earns.

Each citizen wants a short distance between his residence and the office. If the contractor builds the bridge 1/2/3/4, then  $R$  will earn 200/100/150/50 points, and  $B$  will earn 100/200/150/50 points, respectively. Importantly, the contractor cannot simply build any bridge he likes.<sup>13</sup> In order for the contractor to build a bridge, he needs to get permission from at least one of the citizens  $R$  and  $B$ . Table 1 summarizes the points each citizen earns from each bridge built.

---

<sup>13</sup>In every round of the experiment, each subject's screen presented a table summarizing the points from each bridge built but not the corresponding map of the city. For more details, please see the experimental instructions presented in Appendix C.

Bridge Built	$R$	$B$
4 (the longest)	50	50
3 (2nd longest)	150	150
2 (3rd longest)	100	200
1 (the shortest)	200	100

Table 1: Points Earnings from Each Bridge Built

## 4.2 Treatments

Our experimental design involves two treatments: Treatments 2P and 3P. Treatment 2P, our baseline treatment, contains ten 2-player games. All 2-player games share the experimental environment described in the previous section while the number of bridges and the preferences of citizens may differ across games. Treatment 3P contains twenty 3-player games. All 3-player games also share the same environment that is slightly modified to include one more citizen  $G$ . Table 5 in Appendix A presents the ordinal preferences profile of ten games used in Treatment 2P and the ordinal preferences profile of twenty games used in Treatment 3P. Each row refers to a game in Table 5. The players’ preferences are given by some permutation of the integers from 1 to 5. The first integer is the most preferred outcome for the player, and so on.

In Treatment 2P, each participant plays each of the ten 2-player games four times, twice in each role, with the games and role therein appearing in random order. Thus each treatment has 40 rounds (10 games  $\times$  2 roles  $\times$  2 repetitions) in total. The ten 2-player games differ from each other with respect to the players’ ordinal preferences and the number of actions (bridges) available for each player. Among the ten games, six of them are five-action games (Games A1-A6) and the rest are four-action games (Games B1-B4). Table 6 in Appendix A documents what strategy a player would choose for each of the behavioral theories proposed in Definition 2 for the ten 2-player games. There are three games with  $\pi^* = 5$ , two games each with  $\pi^* = 2, 3, 4$ , and one game with  $\pi^* = 1$ . The data from Treatment 2P will allow us to test all but Hypothesis 4.

Treatment 3P is primarily designed to test Hypothesis 4 and serve as a robustness check. It comprises twenty 3-player games created by introducing (two different versions of) a third player to each of the ten two-player games in Treatment 2P. Each participant engages in all twenty 2-player games, assuming each role in a random order. As a result, Treatment 3P consists of a total of 60 rounds (20 games  $\times$  3 roles). The preference of the additional player is deliberately chosen to ensure that the inclusion of the third player leads to meaningful changes in the outcomes predicted by the IDWDS. The theoretical predictions of the twenty 3-player games are presented in Table 8 in Appendix A. Table 7 in Appendix A illustrates the changes in the best prediction  $\pi^*$  resulting from the addition of the third player for each game.

## 4.3 Procedure

We conducted the experiments in English using oTree (Chen, Schonger, and Wickens, 2016) and Zoom in March 2022 at the Hong Kong University of Science and Technology. We

conducted six sessions for each treatment and thus had 12 sessions in total. Each session had 15-21 participants. We recruited 218 participants in total, and each of them participated in only one session. We collected 40 observations per individual in Treatment 2P and 60 observations per individual in Treatment 3P. The experimental session lasted approximately 70 minutes. The average earning was HKD 176.3 ( $\approx$  USD 22.6). The instructions for Treatment 2P can be found in Appendix C.

Upon invitation, participants were instructed to find a quiet location with reliable internet access to remain for the entire duration of the experiment. They joined the designated Zoom meeting using their personal laptop or desktop computer. It was mandatory for all participants to keep their video turned on throughout the experiment. The Zoom settings did not allow for chat communication among participants. To provide the experimental instructions, each participant received an individual oTree link through the Zoom chat message. The instructions were also read aloud to ensure that all participants had access to the same information. A between-subject design was utilized, and a random matching protocol was employed. No feedback was given at the end of each round. Towards the conclusion of the experimental sessions, participants were randomly paired, and one game was selected at random from each pair to calculate their earnings.

## 5 Results

We report experimental results in this section. All results reported are based on the observations from all rounds of decision making because we did not provide any feedback in between rounds, and the overall behavior is stable across rounds. Figures 7 and 8 presented in Appendix B illustrate the stable time-trend in terms of the average earnings and amount of time spent.

### 5.1 Results on aggregate behavior

Figure 2 presents a bar chart that describes the outcomes of three related games from both treatments. It contains three bars: the left bar represents the outcome of a 2-player game (Game A1), and the other two bars represent the outcomes of the two 3-player games (Games A11 and A12) generated by adding a third player to the corresponding 2-player game presented on the left. Figure 9 in Appendix B reports ten bar charts for all games.

Focus on the left bar in each chart describing the outcomes of the 2-player games. Among Game As (with five actions), three games have a substantial proportion of outcomes with  $\bar{x} = 5$  (Games A1, A2, and A4), and they are the only three 2-player games with five actions in which theory predicts that  $\pi^* = \bar{x}$ . Among Game Bs (with four actions), there are only two games (B1 and B2) that have a non-negligible proportion ( $\approx 15\%$ ) of outcomes with  $\bar{x} = 4$ . Even in those two games, however, the modal outcomes are strictly below  $\pi^*$ . The outcomes from the 3-player games with four actions are also consistent with the prediction of  $\bar{x}$  only when  $\pi^* = \bar{x}$ . In all games in which theory predicts that  $\pi^* = \bar{x}$  (Games A11, A12, A21, A22, A31, A41, A42, A51, A61, B11, B12, B22, B31, and B41),  $\bar{x}$  is the modal outcome. Moreover, there is no other game in which  $\bar{x}$  is the modal outcome. This observation provides strong evidence in favor of Hypothesis 1. Thus, we have our first result:



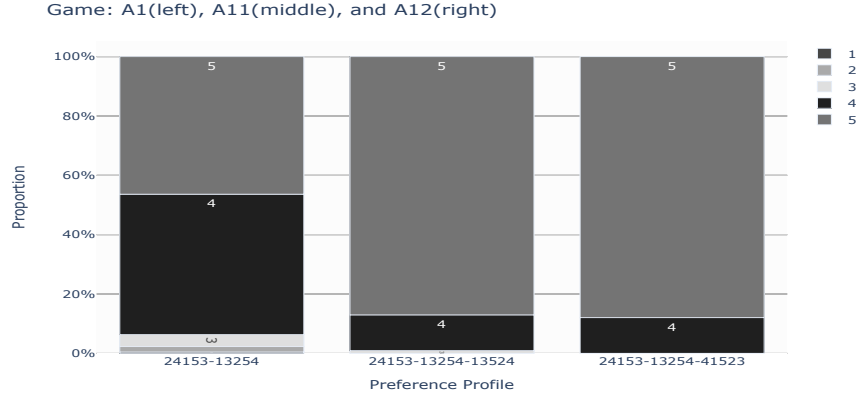


Figure 2: Outcomes – Games A1, A11, and A12

**Result 1.** *In all games with  $\pi^* < \bar{x}$ , the modal outcome of the game is less than  $\bar{x}$ .*

Recall that both Hypotheses 1 and 2 are motivated by Proposition 1. Hypothesis 1 is a weaker statement. The permission game always has a Nash Equilibrium with outcome  $\bar{x}$ . Proposition 1 demonstrates that whenever  $\pi^* < \bar{x}$ , the (unique) outcome of the game that survives iterative deletion of weakly dominated strategies is  $\pi^*$ . Hypothesis 1 states that players will avoid full approval when it is weakly dominated to do so. Hypothesis 2 makes the stronger assertion that players will actually arrive at the outcome that survives IDWDS. Our data confirm this assertion. As presented in Figure 9, the modal outcome observed in every game is  $\pi^*$ .

**Result 2.** *In all games, the modal outcome is  $\pi^*$ .*

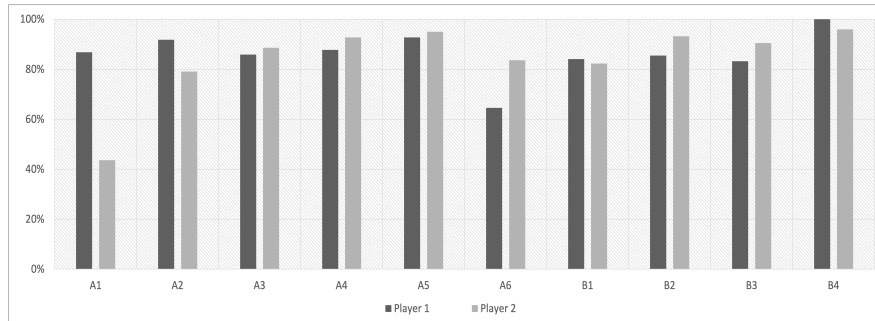


Figure 3: Percentage of Strategies Surviving IDWDS in Treatment 2P

Hypothesis 3, which states that the modal strategy selected survives IDWDS, is stronger than Hypothesis 2 because when the strategy profile is  $\mathbf{x}$ , the outcome of the game is  $M(\mathbf{x})$ . Provided that  $x_i \leq \pi^*$  for all  $i$  and  $x_j = \pi^*$  for some  $j$ ,  $M(\mathbf{x}) = \pi^*$ . There is no need for every player to select a strategy that survives IDWDS. Figure 3 presents the percentage of observations that are consistent with the strategies that survive IDWDS for each game in Treatment 2P. Figure 6 presented in Appendix A reports the percentages for Treatment 3P. Tables 9 and 10 in Appendix A present the percentage of observations (compliance rates)

that are consistent with each of the strategies predicted by all six theories for Treatments 2P and 3P, respectively.<sup>14</sup>

**Result 3.** *For the vast majority of the games, the modal strategies observed survived IDWDS.*

Result 3 underestimates the extent to which players use strategies that survive IDWSD. Taking into account that multiple strategies are consistent with following the maximal iterated deletion of weakly dominated strategy, 85% of observations comply with IDWDS in Treatment 2P and 93% comply in Treatment 3P. They are substantially higher than the IDWSD compliance rates for randomly selected strategies, 36% and 62% for 2P and 3P games, respectively.

Overall, the reported percentage is 85% in the 2P games but there are two exceptions: Player 1 in A6 and Player 2 in A1. In game A6, the player who does not conform is not pivotal. That is, in this game inconsistent individual behavior does not influence the outcome. In A1, it is possible that the players are avoiding the unique equilibrium outcome (5), which is dominated by a non-equilibrium outcome. In Treatment 3P, the observed strategies are more consistent (than in Treatment 2P) with the strategies that survive IDWDS (93%). The higher overall compliance rate is mainly driven by the fact that IDWDS is more permissive in the 3-player games.

We conducted a regression analysis to investigate possible explanations for variations in the compliance rate and Table 2 reports the result for Treatment 2P. Table 2 in Appendix A reports the results for Treatment 3P.

Term	Estimate	SE	p-value	5% CI
(Intercept)	0.759	0.041	0.000	(0.679, 0.839)
Complexity <sub>1</sub>	-0.187	0.016	0.000	(-0.218, -0.156)
Preference <sub>1</sub>	0.038	0.011	0.001	(0.016, 0.061)
Pivotal <sub>1</sub>	0.298	0.030	0.000	(0.238, 0.357)
Complexity <sub>2</sub>	-0.156	0.016	0.000	(-0.187, -0.124)
Preference <sub>2</sub>	0.077	0.011	0.000	(0.055, 0.100)
Pivotal <sub>2</sub>	0.340	0.030	0.000	(0.280, 0.399)

- a. SE refers to the standard error.
- b. CI refers to the confidence interval.
- c. Complexity<sub>*i*</sub> represents the number of local maxima in Player *i*'s preference profile.
- d. Preference<sub>*i*</sub> represents the payoff that Player *i* receives from the best prediction  $\pi^*$ .
- e. Pivotal<sub>*i*</sub> is a binary variable that takes the value 1 if the Player *i*'s decision determines the outcome and 0 otherwise.

Table 2: Linear Regression Analysis for Compliance Rates - Treatment 2P

The analysis identifies three factors that influence compliance. First, if a player's decision determines the outcome (pivotality), the player is more likely to select a strategy consistent with IDWDS. Second, if the BEST prediction  $\pi^*$  yields a higher payoff for a player (preference), the player is more likely to comply. Third, if a player's preference profile has fewer

<sup>14</sup>We decided to report findings on IDWDS when discussing aggregate results and MIDWDS when discussing individual results. MIDWDS and IDWDS differ because MIDWDS makes a unique selection from IDWDS. Because we require theories to be single valued, we must modify IDWDS to make it single valued. MIDWDS is a convenient selection. We discuss IDWDS when describing aggregate results because it is a more fundamental concept. Because any strategy consistent with MIDWDS must be consistent with IDWDS, more subjects comply with IDWDS than MDIWDS.

local maxima (complexity), the player is more likely to comply. These three categories are conceptually orthogonal to each other, and adding interaction terms does not change the qualitative results. Moreover, the results are consistent between Treatments 2P and 3P.

Hypothesis 4 conjectured that adding a player would increase the outcomes based on Table 7 in Appendix A. The data confirm this hypothesis. The bar charts presented in Figure 9 in Appendix B confirm that adding the third player strictly increases either the modal outcome itself (e.g., in the comparison between Game A3 and A31) or the percentage of the same modal outcome (e.g., in the comparison between Game A3 and A32) for all games.

**Result 4.** *Adding a player increases the outcome.*

## 5.2 Results on individual behavior

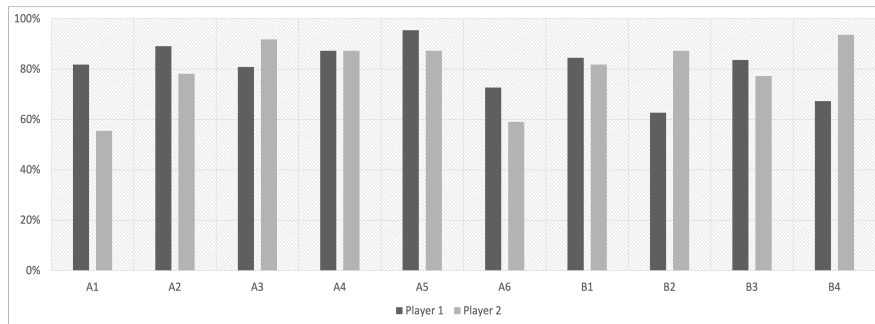


Figure 4: Consistency in Treatment 2P

Figure 4 reports the consistency rate aggregated across all individuals for each game in Treatment 2P. Recall that every subject in Treatment 2P played each game two times in each role. We say that his/her choice is consistent in a game if the same strategy is chosen in identical games, and inconsistent otherwise. The overall consistency rate is slightly above 80%, but there are four cases with which the consistency rate is substantially lower than the average: Player 2 in A1 and A6, and Player 1 in B2 and B4. In the latter three cases, players who are categorized as L1 and L2 type are indifferent between multiple actions. As we suggested in our discussion of Result 2, inconsistency in A6 and B2 may be due to a player recognizing that her action is not pivotal.

**Result 5.** *Subjects play the same strategy in identical games with probability 0.8.*

Consistency requires that when a player faces the same decision problem, she behaves in the same way. The best prediction we can make is that a subject will repeat her earlier choice the second time she plays a game. Consistency is far from perfect, however. Subjects play different strategies in identical games 20% of the time. Failure to be consistent could be due to errors, indifference, boredom, experimentation, learning, hedging, or a desire to randomize.

Result 5 is a minimal consistency condition. We would like to be able to use observations of behavior to predict behavior in other situations. Predictable players will play the same game in the same way. Of course, the hypothesis assumes that subjects treat a particular payoff matrix in the same way no matter where in the sequence of games it appears.

Now we report the results from our cross-validation exercise. Table 12 in Appendix B presents the test and training scores for each possible combination of the six behavioral theories. The scores reported are the averages from 100 randomly selected subsets of size 30 for the training set and size 10 for the test set in the 2P games, and size 45 for the training set and size 15 for the test set in the 3P games.

**Result 6.** *Test scores are below consistency scores.*

Result 6 corresponds to Hypothesis 6. Consistency scores are indeed higher than test scores, but the difference is quite small. This suggests that we can attribute most of our failure to predict behavior to baseline randomness. Subjects fail to play according to the rule that best describes their behavior at essentially the same rate that they fail to play identical games in the same way.

Although test scores are below consistency scores, they are not much below. Players use the same strategy in the same game roughly 80% of the time. Players play according to the category assigned to them roughly 79% of the time. It is plausible to view 80% as an upper bound of the accuracy of our predictions. Predictions based on our classifications approximate this upper bound, but we cannot classify behavior without using a training set, which permits us to make different predictions for different subjects. The predictions of consistency hypothesis hold uniformly for all subjects (and therefore do not require observation of behavior on a training sample).

It is natural to try to relate consistency with other behavior. In two-player games we can rank players by how well they are classified – the number of instances in which their decisions coincide with the prediction of the theory that best describes their behavior (strong consistency). We can also rank players by their consistency – the number of instances in which they make the same choice in identical situations (weak consistency). We would expect the two notions of consistency to be (at least) weakly positively correlated because all of our theories predict identical behavior in identical games. Figure 10 in Appendix B presents the correlation between them. The Spearman coefficient is 0.457 (moderate-to-strong degree of correlation).<sup>15</sup>

**Result 7.** *Different agents are best described by different behavioral rules.*

Size	Best Subset	Test Score	Training Score
1	L2-p	0.7824	0.7674
2	L0, L2-p	0.7866	0.7839
3	L0, L2-f, L2-p	0.7915	0.7942
4	L0, L1-f, L2-f, L2-p	0.7919	0.7987
5	L0, L1-f, L2-f, L1-p, L2-p	0.7890	0.8006
6	L0, L1-f, L2-f, L1-p, L2-p, MIDWD	0.7838	0.8082

Table 3: Test and Training Scores form the Best Subset of Each Size - Treatments 2P & 3P

We expect training scores to be higher than test scores because we select theories to do well on the training set. This property does not hold uniformly: Test scores are higher when

<sup>15</sup>The Spearman coefficient is lowered to 0.380 (moderate to strong correlation) when we do count behavior in games that are repeated only once. This difference comes from the fact that all theories imply consistency.

there are only one or two theories. We attribute this to elements in training sets as being more difficult than elements in test sets. We can confirm (both by computation and theory) that if one compares training and test scores averaged over randomly selected sets, then training scores are always greater than test scores.

Result 7 provides evidence in favor of our Hypothesis 7. Table 3 presents the test and training scores from the best subset of each size where data from both treatments is used for the classification. Tables 13 and 14 in Appendix B provide the test scores calculated based on each treatment separately. We obtain higher classification scores when there is more than one theory. The best single theory is  $L2-p$ , which has a test score of .7824. This score increases to .7919 when we can classify agents using  $L0$ ,  $L1-f$ , and  $L2-p$  as well. When we limit attention to only three-player games with five strategies (see Table 18 in Appendix B), the highest test score comes from the one-theory set of  $L2-p$ . For this fraction of the sample, Hypothesis 7 fails.

**Result 8.** *Scores on the training sets are higher than scores on the test set.*

Result 8 confirms Hypothesis 8. The best family of theories predicts correctly 79.19% of the time. The corresponding score on the training set is 79.87%. The (one-sided) paired-sample Wilcoxon test reveals that we can reject the null hypothesis that these two values are the same, in favor of the alternative that the difference is significant ( $p$ -value = 0.0898).<sup>16</sup> We expect the training score to be higher because we select the theories that best fit the training set. The fact that our test scores are close to the training scores suggests that we have approximated the upper bound on how well we can organize the data.

**Result 9.** *Scores are higher in two-player games than in three-player games.*

Tables 13 and 14 in Appendix B show that the scores on three-player games are about 3 percent lower than the test scores on two-player games. Result 9 suggests that two-player games are simpler than three-player games in the sense that we are able to make better classifications.

**Result 10.** *Scores in four-strategy games are essentially identical to scores in five-strategy games.*

Tables 15, 16, 17, and 18 in Appendix B report test scores separately for four-strategy games and five-strategy games. In contrast to Hypothesis 10, Result 10 suggests that four-strategy games are no simpler than five-strategy games in the sense that classification results are nearly identical (approximately 79% in each case). On one hand, this result is encouraging, because it suggests that our theories may help organize observations in games with larger strategy sets. On the other hand, we are skeptical that we will predict behavior as accurately in games with hundreds of strategies.

**Result 11.** *The set of theories that maximize test scores usually does not include MIDWD.*

---

<sup>16</sup>For this analysis, we compare the training score from the best subset and the corresponding test score for each individual. Then using the individual scores as independent observations, we conduct the paired-sample Wilcoxon test.

Result 11 demonstrates that Hypothesis 11 does not hold in general.<sup>17</sup> Iterated deletion of weakly dominated strategies makes good predictions in many cases, but variations of level-2 behavior, which agree with the weak dominance prediction frequently, always appear as one of the best theories, while MIDWD does not. Observe that MIDWD removes weakly dominated strategies using a particular order. In general, more strategies are consistent with the removal of weakly dominated strategies (78.2% in Treatment 2P and 83.2% in Treatment 3P).<sup>18</sup>

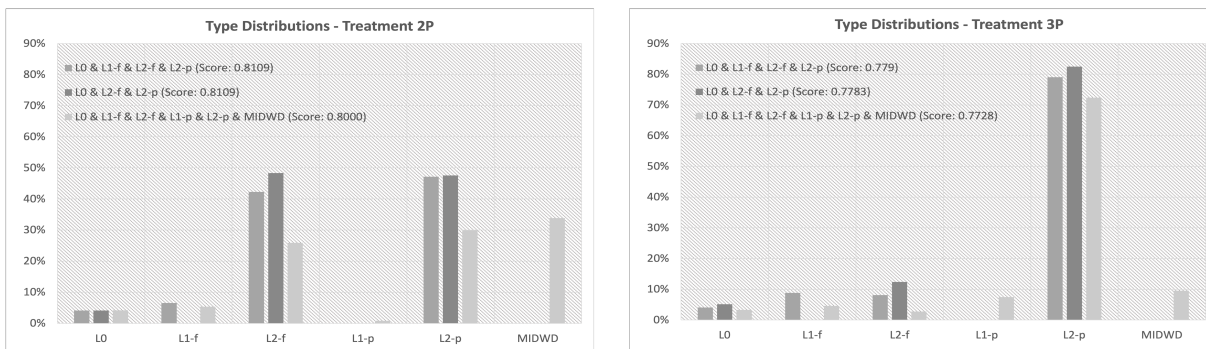


Figure 5: Type Distributions

Figure 5 presents the distribution of behavioral types based on the first- and second-best subsets, as well as the highest-scored subset that involves MIDWD.<sup>19</sup> One feature of our characterization results is that the test score is maximized using four theories. In both two- and three-player games adding a fifth or sixth theory (slightly) reduces the quality of predictions. This finding indicates that using too many theories could result in overfitting. Adding a theory typically increases the score on the test sample,<sup>20</sup> but it could lead to selecting a theory that works well for the training set but makes poor out-of-sample predictions.

Another feature of our characterization is the family of rules that makes good predictions. Level 0 behavior is part of the set of theories that best identify behavior. Not many agents behave according to level 0, but there are 4-5% of individuals who are consistently categorized as level 0 regardless of whether we use the first-, second-, or third-best subset. Note that on its own, level 0 performs badly (the test score is 0.477).

### 5.3 Robustness

We conducted several robustness exercises.

We conducted our principal analysis using the first 30 and 45 games as a training set and the final 10 and 15 games as a test set for Treatments 2P and 3P, respectively. We checked to see whether these choices influenced the results.

<sup>17</sup>It does hold for two-player games and three theories and three-player games with four or five theories.

<sup>18</sup>One can see from Table 8 that  $L2-p$  specifies a weakly higher strategy than MIDWD for all players in Treatment 3P and that when MIDWD and  $L2-p$  differ for a player, the outcome does not depend on whether the player follows MIDWD or  $L2-p$ .

<sup>19</sup>Figure 11 reports the distribution in which both data from Treatments 2P and 3P are used for the classifications.

<sup>20</sup>Either the new theory is better for a subject and training score goes up or it doesn't influence score.



Treatment 2P	30/10 Split	29/10 Split	28/10 Split
	0.8109	0.8109	0.8141
Treatment 3P	45/15 Split	44/15 Split	43/15 Split
	0.7815	0.7796	0.7790

Table 4: Robustness of Best Test Scores

We varied the number of training rounds from 28 to 30 for Treatment 2P and 43 to 45 for Treatment 3P, and the result is reported in Table 4. Doing so changes the test scores slightly (aggregate and for subsamples), but never by more than .4 percent. In some cases, the best collection of theories changes, but the performance of the collection that we originally identified as best is never more than .2 percent lower than the alternative.

Figure 12 illustrates the impact of increasing the size of the training set. It leads to small improvements in test scores until the training set reaches about 30, after which the scores level off. In three-player games, the best prediction shows no improvement in test scores even when the training set is increased from a low number because the single-theory  $L2-p$  approximates the best prediction. If we conduct the increase in the training set size by comparing test scores using the first  $k$  games to train and varying  $k$ , the test scores are not necessarily monotonically increasing (although the score for  $k = 1$  is lower than  $k = 30$ ). Randomizing the membership of the training set smoothes the data. Average test scores then increase with the cardinality of the training set.

It is possible that using two-player (or four-strategy) games in the training set will lead to different predictions in two-player (four-strategy) games than in three-player (five-strategy) games. We investigated whether the results of our classification exercise depends on the number of players or strategies in either the training or test sets. We found no relationship between the quality of test scores as a function of characteristics of either training or test sets. That is, for example, we do not obtain worse predictions on five-strategy games if we train only on four-strategy games. This finding suggests that we can reliably extrapolate play in four-strategy games to forecast behavior in larger games.

## References

- Agranov, Marina and Pietro Ortoleva. 2017. “Stochastic choice and preferences for randomization.” *Journal of Political Economy* 125 (1):40–68.
- Arad, Ayala and Ariel Rubinstein. 2012. “The 11–20 money request game: A level-k reasoning study.” *American Economic Review* 102 (7):3561–3573.
- Battaglini, Marco. 2002. “Multiple referrals and multidimensional cheap talk.” *Econometrica* 70 (4):1379–1401.
- Camerer, Colin F. 2003. “Behavioural studies of strategic thinking in games.” *Trends in cognitive sciences* 7 (5):225–231.
- Chen, Daniel L., Martin Schonger, and Chris Wickens. 2016. “oTree—An open-source platform for laboratory, online, and field experiments.” *Journal of Behavioral and Experimental Finance* 9:88–97.
- Crawford, Vincent P. 2003. “Lying for strategic advantage: Rational and boundedly rational misrepresentation of intentions.” *American Economic Review* 93 (1):133–149.
- Crawford, Vincent P, Miguel A Costa-Gomes, and Nagore Iriberri. 2013. “Structural models of nonequilibrium strategic thinking: Theory, evidence, and applications.” *Journal of Economic Literature* 51 (1):5–62.
- Hastie, Trevor, Robert Tibshirani, and Jerome H Friedman. 2009. *The elements of statistical learning: data mining, inference, and prediction*, vol. 2. Springer.
- Hu, Peicong and Joel Sobel. 2022. “Getting permission.” *American Economic Review: Insights* 4 (4):459–472.
- . 2023. “Getting permission when options are partially ordered.” In *David Gale: Mathematical Economist: Essays in Appreciation on his 100<sup>th</sup> birthday*, edited by M. Ali Khan, Nobusumi Sagara, and Alexander Zaslavski. Springer.
- Kohlberg, Elon and Jean-Francois Mertens. 1986. “On the Strategic Stability of Equilibria.” *Econometrica* 54 (5):1003–1037. URL <http://www.jstor.org/stable/1912320>.
- Krishna, Vijay and John Morgan. 2001. “A model of expertise.” *The Quarterly Journal of Economics* 116 (2):747–775.
- Lai, Ernest K, Wooyoung Lim, and Joseph Tao-yi Wang. 2015. “An experimental analysis of multidimensional cheap talk.” *Games and Economic Behavior* 91:114–144.
- Nagel, Rosemarie. 1995. “Unraveling in Guessing Games: An Experimental Study.” *American Economic Review* 85 (5):1313–26. URL <http://ideas.repec.org/a/aea/aecrev/v85y1995i5p1313-26.html>.
- Nielsen, Kirby and John Rehbeck. 2022. “When choices are mistakes.” *American Economic Review* 112 (7):2237–2268.

- Stahl, Dale O. and Paul W. Wilson. 1994. “Experimental evidence on players’ models of other players.” *Journal of Economic Behavior and Organization* 25 (3):309 – 327. URL <http://www.sciencedirect.com/science/article/pii/0167268194901031>.
- . 1995. “On players’ models of other players: theory and experimental evidence.” *Games and Economic Behavior* 10 (1):218 – 254. URL <http://www.sciencedirect.com/science/article/pii/S0899825685710317>.
- Stone, Mervyn. 1974. “Cross-validators choice and assessment of statistical predictions.” *Journal of the royal statistical society: Series B (Methodological)* 36 (2):111–133.
- Vespa, Emanuel and Alistair J Wilson. 2016. “Communication with multiple senders: An experiment.” *Quantitative Economics* 7 (1):1–36.

# Appendices

## A Prediction and Compliance Rates

Table 5: Experimental Games and Ordinal Preferences

Treatment 2P			Treatment 3P			
Game	P1	P2	Game	P1	P2	P3
A1	24153	13254	A11	24153	13254	13524
			A12	24153	13254	41523
A2	24153	31524	A21	24153	31524	13524
			A22	24153	31524	41523
A3	31245	14253	A31	31245	14253	13524
			A32	31245	14253	41523
A4	13524	42315	A41	13524	42315	13524
			A42	13524	42315	41523
A5	24153	31245	A51	24153	31245	13524
			A52	24153	31245	41523
A6	24153	12345	A61	24153	12345	13524
			A62	24153	12345	41523
B1	1243	1324	B11	1243	1324	2431
			B12	1243	1324	3142
B2	1243	2134	B21	1243	2134	2431
			B22	1243	2134	3142
B3	1324	2134	B31	1324	2134	2431
			B32	1324	2134	3142
B4	1324	3142	B41	1324	3142	2431
			B42	1324	3142	3142

Table 6: Predictions - Treatment 2P

Game	P1's Strategy							P2's Strategy							Outcome
	$L_0$	$L_1^f$	$L_2^f$	$L_1^p$	$L_2^p$	MID	IDWDS	$L_0$	$L_1^f$	$L_2^f$	$L_1^p$	$L_2^p$	MID	IDWDS	$\pi^*$
A1	2	2	4	2	4	4	4,5	1	3	3	3	3	5	5	5
A2	2	4	4	4	4	4	4,5	3	3	5	3	5	5	5	5
A3	3	3	3	3	3	3	2,3,4	1	4	4	4	4	4	4	4
A4	1	5	5	5	5	5	5	4	4	4	4	4	4	2,4,5	5
A5	2	4	4	4	4	4	4	3	3	3	3	3	3	1,3,4	4
A6	2	2	2	2	2	2	2	1	1	1	1	1	1	1,2	2
B1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
B2	1	1	1	1	1	2	1,2	2	2	2	2	2	2	2	2
B3	1	3	3	3	3	3	3	2	2	2	2	2	2	2,3	3
B4	1	1	1	3	3	1	1,2,3	3	3	3	3	3	3	3	3

- <sup>a.</sup>  $P_i$  refers to the Player  $i$  for  $i = 1, 2$ .
- <sup>b.</sup>  $L_0$  refers to the favorite action.  $L_i$  with  $i > 0$  refers to the best response to the opponent  $L_{(i-1)}$ .
- <sup>c.</sup> MID refers to the strategy that survives the maximal iterated deletion of weakly dominated (MIDWD) strategies.
- <sup>d.</sup> IDWDS refers to the strategy that survives any order of iterated deletion of weakly dominated strategies.
- <sup>e.</sup>  $\pi^*$  refers to the outcome from the best equilibrium for both players.

Table 7: Does adding Player 3 change the best outcome  $\pi^*$ ?

	Ordinal Preference		Player 3's Preference			
	Player 1	Player 2	13524	41523	2431	3142
A1	24153	13254	No	No	N/A	
A2	24153	31524	No	No		
A3	31245	14253	4 to 5	No		
A4	13524	42315	No	No		
A5	24153	31245	4 to 5	No		
A6	24153	12345	2 to 5	2 to 4		
B1	1243	1324	N/A		1 to 4	1 to 4
B2	1243	2134			No	2 to 4
B3	1324	2134			3 to 4	No
B4	1324	3142			3 to 4	No

- <sup>a.</sup> The entry “No” implies that adding the third player does not change the best equilibrium prediction from its base two-player game.
- <sup>a.</sup> The entry “X to Y” implies that adding the third player changes the best equilibrium prediction from X to Y.

Table 8: Predictions - Treatment 3P

Game	P1's Strategy						P2's Strategy						P3's Strategy						Outcome $\pi^*$			
	$L_0$	$L_1^f$	$L_2^f$	$L_1^p$	$L_2^p$	MID	IDWDS	$L_0$	$L_1^f$	$L_2^f$	$L_1^p$	$L_2^p$	MID	IDWDS	$L_0$	$L_1^f$	$L_2^f$	$L_1^p$		$L_2^p$	MID	IDWDS
A11	2	2	4	2	4	4	4	1	3	1	3	3	5	4,5	1	3	1	3	3	5	5	4,5
A12	2	2	2	4	4	2	1,2,3,4,5	1	5	5	5	5	5	4,5	4	4	4	4	4	4	4	4,5
A21	2	4	4	4	4	4	1,2,4	3	3	5	3	5	5	1,2,3,4,5	1	1	5	3	5	5	5	1,2,3,4,5
A22	2	2	2	4	4	2	1,2,3,4,5	3	5	5	5	5	5	5	4	4	4	4	4	4	4	1,2,3,4,5
A31	3	3	3	3	3	3	2,3,4,5	1	4	4	4	4	4	4,5	1	1	5	3	5	5	5	5
A32	3	3	3	3	3	3	1,2,3,4	1	1	1	4	4	4	1,2,3,4	4	4	4	4	4	4	4	1,2,3,4
A41	1	5	1	5	5	5	1,2,3,4,5	4	4	4	4	4	4	4	1	5	1	5	5	5	5	1,2,3,4,5
A42	1	5	5	5	5	5	5	4	4	4	4	4	4	1,2,3,4	4	4	4	4	4	4	4	1,2,3,4
A51	2	4	4	4	4	4	1,2,4,5	3	3	3	3	3	3	2,3,4,5	1	1	5	3	5	5	5	5
A52	2	2	2	4	4	2	2,4	3	3	3	3	3	3	1,2,3,4	4	4	4	4	4	4	4	4
A61	2	2	4	2	4	4	4	1	1	1	1	1	1	1,2,3,4	1	3	3	3	3	5	5	5
A62	2	2	2	4	4	2	1,2,3,4	1	1	1	1	1	1	1,2,3,4	4	4	4	4	4	4	4	3,4
B11	1	1	4	1	4	4	1,2,3,4	1	3	3	3	3	3	3,4	2	2	4	2	4	4	4	2,3,4
B12	1	4	4	4	4	4	4	1	1	1	3	3	1	1,2,3,4	3	3	3	3	3	3	3	2,3,4
B21	1	1	1	1	1	1	1,2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	1,2
B22	1	4	4	4	4	4	4	2	2	2	2	2	2	1,2,3,4	3	3	3	3	3	3	3	3,4
B31	1	3	3	3	3	3	1,3,4	2	2	2	2	2	2	1,2,3,4	2	2	4	2	4	4	4	4
B32	1	1	1	3	3	1	1,2,3	2	2	2	2	2	2	1,2,3	3	3	3	3	3	3	3	2,3
B41	1	1	1	3	3	1	1,2,3,4	3	3	3	3	3	3	1,3,4	2	4	4	4	4	4	4	4
B42	1	1	1	3	3	1	1,2,3	3	3	3	3	3	3	1,2,3	3	3	3	3	3	3	3	1,2,3

a.  $P_i$  refers to the Player  $i$  for  $i = 1, 2$ .

b.  $L_0$  refers to the favorite action.  $L_i$  with  $i > 0$  refers to the best response to the opponent  $L_{(i-1)}$ .

c. MID refers to the strategy that survives the maximal iterated deletion of weakly dominated (MIDWD) strategies.

d. IDWDS refers to the strategy that survives any order of iterated deletion of weakly dominated strategies.

e.  $\pi^*$  refers to the outcome from the best equilibrium for both players.



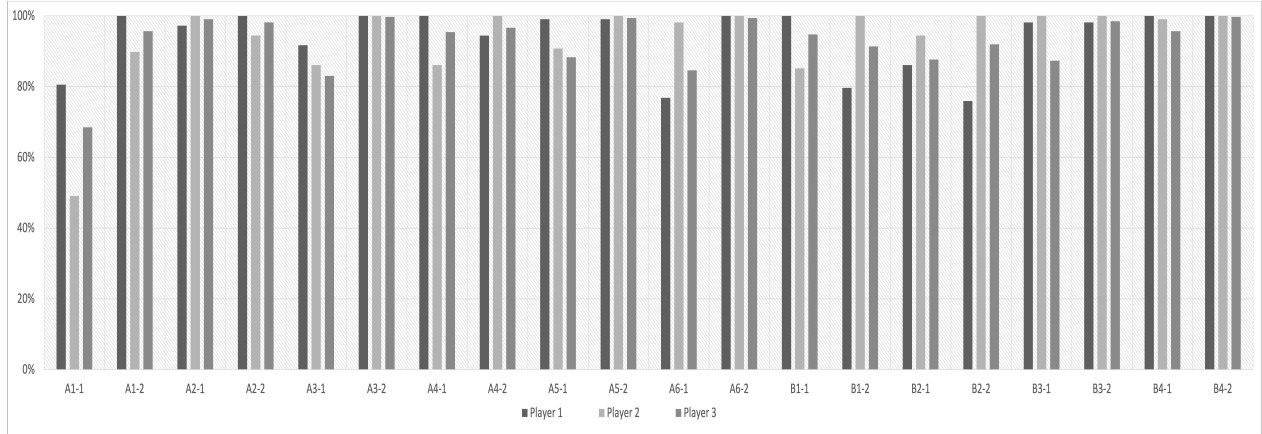


Figure 6: Percentage of Strategies Surviving IDWDS in Treatment 3P

Table 9: Compliance Rates - Treatment 2P

Game	P1's Strategy							P2's Strategy							Outcome
	L0	L1-f	L2-f	L1-p	L2-p	MID	IDWDS	L0	L1-f	L2-f	L1-p	L2-p	MID	IDWDS	$\pi^*$
A1	0.11	0.11	0.82	0.11	0.82	0.82	0.87	0.12	0.35	0.35	0.35	0.35	0.44	0.44	0.46
A2	0.05	0.86	0.86	0.86	0.86	0.86	0.92	0.16	0.16	0.79	0.16	0.79	0.79	0.79	0.80
A3	0.76	0.76	0.76	0.76	0.76	0.76	0.86	0.10	0.89	0.89	0.89	0.89	0.89	0.89	0.90
A4	0.03	0.88	0.88	0.88	0.88	0.88	0.88	0.88	0.88	0.88	0.88	0.88	0.88	0.93	0.88
A5	0.05	0.93	0.93	0.93	0.93	0.93	0.93	0.84	0.84	0.84	0.84	0.84	0.84	0.95	0.93
A6	0.65	0.65	0.65	0.65	0.65	0.65	0.65	0.40	0.40	0.40	0.40	0.40	0.40	0.40	0.55
B1	0.84	0.84	0.84	0.84	0.84	0.84	0.84	0.82	0.82	0.82	0.82	0.82	0.82	0.82	0.69
B2	0.55	0.55	0.55	0.55	0.55	0.30	0.85	0.93	0.93	0.93	0.93	0.93	0.93	0.93	0.80
B3	0.11	0.83	0.83	0.83	0.83	0.83	0.83	0.83	0.83	0.83	0.83	0.83	0.83	0.90	0.82
B4	0.54	0.54	0.54	0.45	0.45	0.54	1.00	0.96	0.96	0.96	0.96	0.96	0.96	0.96	0.97
Average	0.37	0.70	0.77	0.69	0.76	0.74	0.86	0.60	0.17	0.77	0.71	0.77	0.78	0.84	0.78

<sup>a.</sup>  $P_i$  refers to the Player  $i$  for  $i = 1, 2$ .

<sup>b.</sup>  $L_0$  refers to the favorite action.  $L_i$  with  $i > 0$  refers to the best response to the opponent  $L_{(i-1)}$ .

<sup>c.</sup> MID refers to the strategy that survives the maximal iterated deletion of weakly dominated (MIDWD) strategies.

<sup>d.</sup> IDWDS refers to the strategy that survives any order of iterated deletion of weakly dominated strategies.

<sup>e.</sup>  $\pi^*$  refers to the outcome from the best equilibrium for both players.

Table 10: Compliance Rate - Treatment 3P

Game	P1's Strategy						P2's Strategy						P3's Strategy						Outcome			
	$L_0$	$L_1^f$	$L_2^f$	$L_1^p$	$L_2^p$	MID	IDWDS	$L_0$	$L_1^f$	$L_2^f$	$L_1^p$	$L_2^p$	MID	IDWDS	$L_0$	$L_1^f$	$L_2^f$	$L_1^p$		$L_2^p$	MID	IDWDS
A1-1	0.15	0.15	0.81	0.15	0.81	0.81	0.81	0.15	0.36	0.15	0.36	0.36	0.46	0.49	0.06	0.17	0.06	0.17	0.17	0.75	0.76	0.87
A1-2	0.19	0.19	0.19	0.78	0.78	0.19	1.00	0.06	0.88	0.88	0.88	0.88	0.88	0.90	0.96	0.96	0.96	0.96	0.96	0.96	0.97	0.88
A2-1	0.06	0.90	0.90	0.90	0.90	0.90	0.97	0.21	0.21	0.78	0.21	0.78	0.78	1.00	0.11	0.11	0.78	0.10	0.78	0.78	1.00	0.93
A2-2	0.16	0.16	0.16	0.80	0.80	0.16	1.00	0.03	0.94	0.94	0.94	0.94	0.94	0.94	0.96	0.96	0.96	0.96	0.96	0.96	1.00	0.94
A3-1	0.79	0.79	0.79	0.79	0.79	0.79	0.92	0.13	0.85	0.85	0.85	0.85	0.85	0.86	0.20	0.20	0.71	0.08	0.71	0.71	0.71	0.71
A3-2	0.69	0.69	0.69	0.69	0.69	0.69	1.00	0.19	0.19	0.19	0.78	0.78	0.78	1.00	0.97	0.97	0.97	0.97	0.97	0.97	0.99	0.99
A4-1	0.05	0.83	0.05	0.83	0.83	0.83	1.00	0.86	0.86	0.86	0.86	0.86	0.86	0.86	0.04	0.88	0.04	0.88	0.88	0.88	1.00	0.97
A4-2	0.05	0.94	0.94	0.94	0.94	0.94	0.94	0.97	0.97	0.97	0.97	0.97	0.97	1.00	0.93	0.93	0.93	0.93	0.93	0.93	0.95	0.95
A5-1	0.06	0.91	0.91	0.91	0.91	0.91	0.99	0.82	0.82	0.82	0.82	0.82	0.82	0.91	0.10	0.10	0.75	0.14	0.75	0.75	0.75	0.75
A5-2	0.14	0.14	0.14	0.85	0.85	0.14	0.99	0.74	0.74	0.74	0.74	0.74	0.74	1.00	0.99	0.99	0.99	0.99	0.99	0.99	0.99	1.00
A6-1	0.17	0.17	0.77	0.17	0.77	0.77	0.77	0.50	0.50	0.50	0.50	0.50	0.50	0.98	0.06	0.14	0.14	0.14	0.14	0.79	0.79	0.80
A6-2	0.22	0.22	0.22	0.77	0.77	0.22	1.00	0.48	0.48	0.48	0.48	0.48	0.48	1.00	0.98	0.98	0.98	0.98	0.98	0.98	0.98	0.98
B1-1	0.26	0.26	0.48	0.26	0.48	0.48	1.00	0.12	0.83	0.83	0.83	0.83	0.83	0.85	0.40	0.40	0.59	0.40	0.59	0.59	0.99	0.79
B1-2	0.18	0.80	0.80	0.80	0.80	0.80	0.80	0.42	0.42	0.42	0.57	0.42	1.00	0.93	0.93	0.93	0.93	0.93	0.93	0.94	0.94	0.81
B2-1	0.28	0.28	0.28	0.28	0.28	0.28	0.86	0.94	0.94	0.94	0.94	0.94	0.94	0.94	0.82	0.82	0.82	0.82	0.82	0.82	0.82	0.69
B2-2	0.22	0.76	0.76	0.76	0.76	0.76	0.76	0.64	0.64	0.64	0.64	0.64	0.64	0.64	0.94	0.94	0.94	0.94	0.94	1.00	0.78	
B3-1	0.10	0.88	0.88	0.88	0.88	0.88	0.98	0.89	0.89	0.89	0.89	0.89	0.89	0.89	0.34	0.34	0.64	0.34	0.64	0.64	0.64	0.65
B3-2	0.30	0.30	0.30	0.67	0.67	0.30	0.98	0.69	0.69	0.69	0.69	0.69	0.69	0.69	0.97	0.97	0.97	0.97	0.97	0.97	0.97	0.96
B4-1	0.19	0.19	0.19	0.81	0.81	0.19	1.00	0.96	0.96	0.96	0.96	0.96	0.96	0.96	0.05	0.88	0.88	0.88	0.88	0.88	0.88	0.88
B4-2	0.35	0.35	0.35	0.64	0.64	0.35	1.00	0.98	0.98	0.98	0.98	0.98	0.98	0.98	0.98	0.98	0.98	0.98	0.98	0.98	0.99	0.99
Average	0.23	0.50	0.53	0.68	0.76	0.57	0.94	0.54	0.71	0.73	0.75	0.77	0.77	0.94	0.59	0.68	0.75	0.68	0.80	0.86	0.91	0.87

a.  $P_i$  refers to the Player  $i$  for  $i = 1, 2$ .

b.  $L_0$  refers to the favorite action.  $L_i$  with  $i > 0$  refers to the best response to the best response to the best response to the opponent  $L_{(i-1)}$ .

c. MID refers to the strategy that survives the maximal iterated deletion of weakly dominated (MIDWD) strategies.

d. IDWDS refers to the strategy that survives any order of iterated deletion of weakly dominated strategies.

e.  $\pi^*$  refers to the outcome from the best equilibrium for both players.

## B Additional Figures and Tables

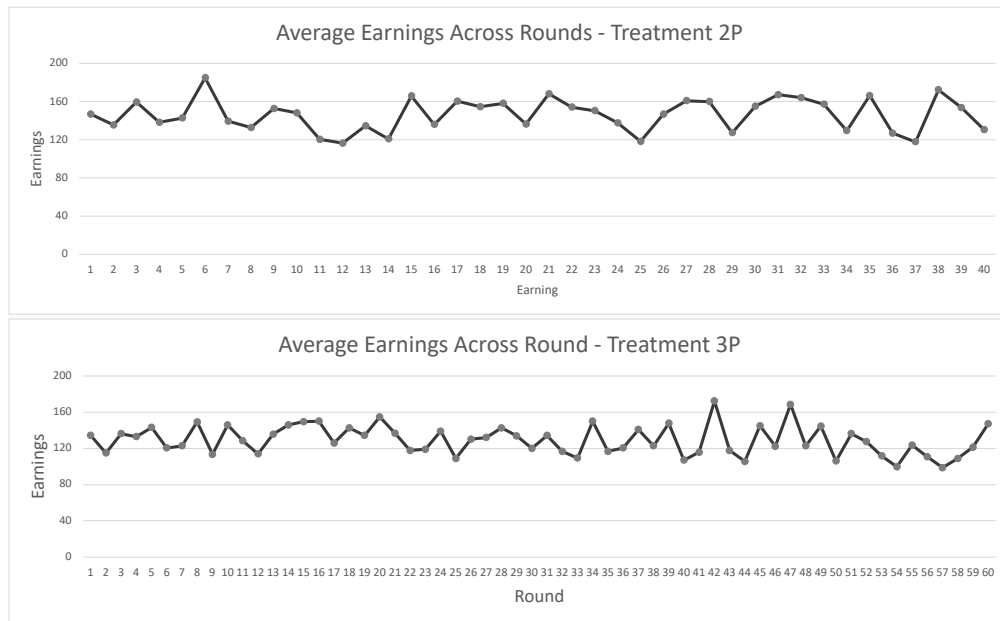


Figure 7: Average Earning

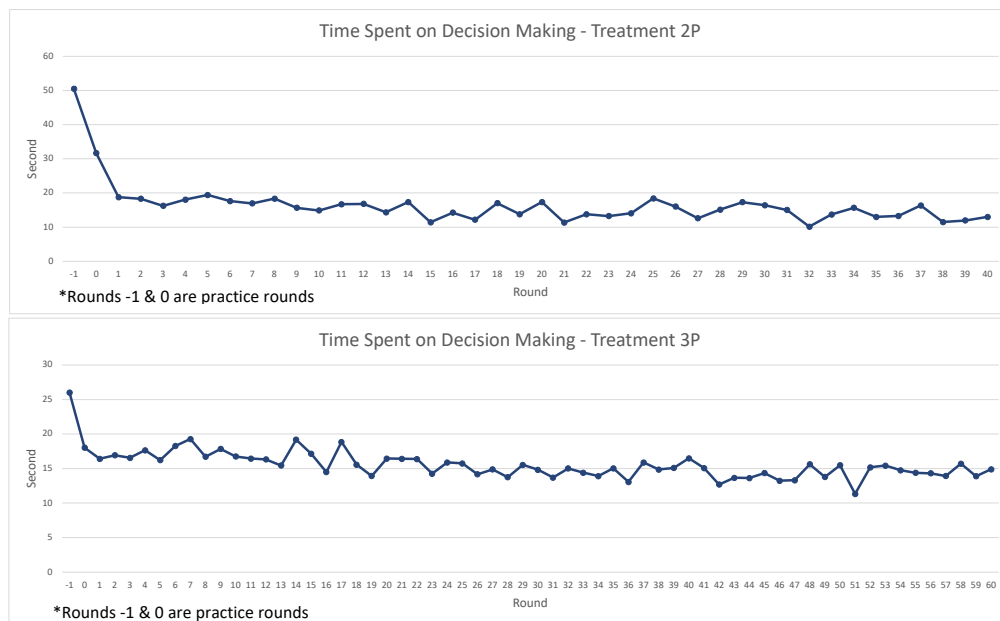


Figure 8: Time Spent In Each Round

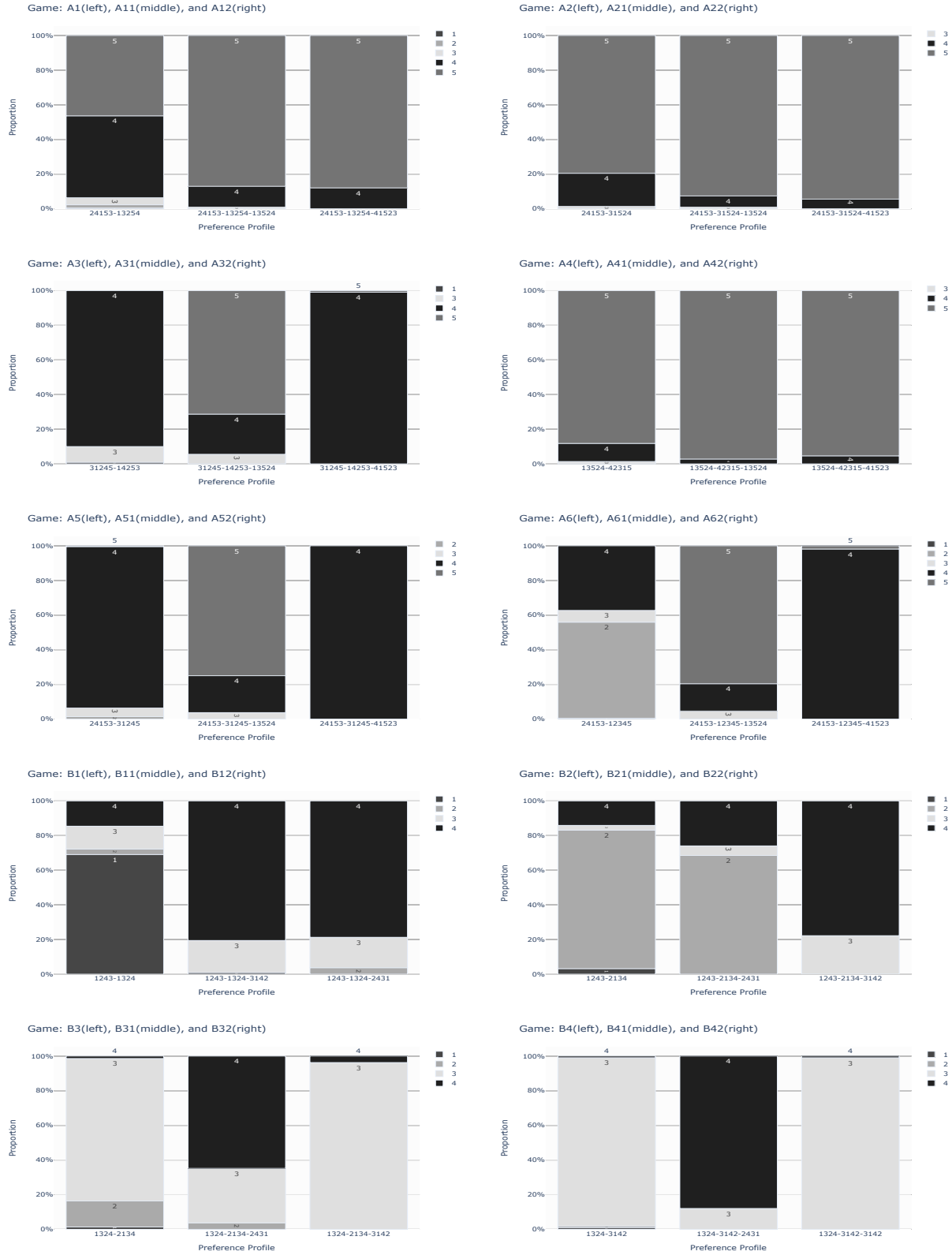


Figure 9: Outcomes (All Games)

Term	Estimate	SE	p-value	5% CI
(Intercept)	0.779	0.031	0.000	(0.719, 0.839)
Complexity <sub>1</sub>	-0.118	0.011	0.000	(-0.139, -0.097)
Preference <sub>1</sub>	0.036	0.006	0.000	(0.025, 0.047)
Pivotal <sub>1</sub>	0.071	0.013	0.000	(0.045, 0.097)
Complexity <sub>2</sub>	-0.069	0.011	0.000	(-0.090, -0.048)
Preference <sub>2</sub>	0.048	0.006	0.000	(0.037, 0.059)
Pivotal <sub>2</sub>	0.077	0.013	0.000	(0.051, 0.103)
Complexity <sub>3</sub>	-0.088	0.011	0.000	(-0.109, -0.067)
Preference <sub>3</sub>	0.042	0.006	0.000	(0.031, 0.053)
Pivotal <sub>3</sub>	0.070	0.013	0.000	(0.044, 0.096)

- a. SE refers to the standard error.
- b. CI refers to the confidence interval.
- c. Complexity<sub>*i*</sub> represents the number of local maxima in Player *i*'s preference profile.
- d. Preference<sub>*i*</sub> represents the payoff that Player *i* receives from the best prediction  $\pi^*$ .
- e. Pivotal<sub>*i*</sub> is a binary variable that takes the value 1 if the Player *i*'s decision determines the outcome and 0 otherwise.

Table 11: Linear Regression Analysis for Compliance Rates - Treatment 3P

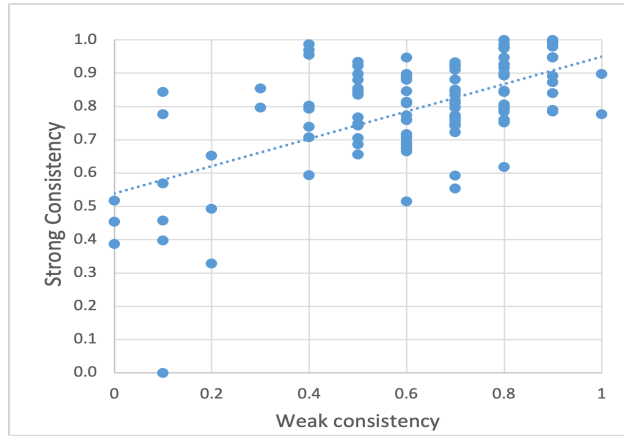


Figure 10: Correlation Between Strong and Weak Consistencies

Table 12: Test and Training Scores for All Subsets of Behavioral Theories

Size	Subset	Test Scores	Training Scores	Size	Subset	Test Scores	Training Scores
1	L0	0.4648	0.4671	2	L0, L1-f	0.6585	0.6674
	L1-f	0.6523	0.6596		L0, L2-f	0.7299	0.7374
	L2-f	0.7212	0.7264		L0, L1-p	0.7040	0.7138
	L1-p	0.6944	0.7016		L0, L2-p	0.7805	0.7870
	L2-p	0.7679	0.7722		L0, MIDWD	0.7502	0.7564
	MIDWD	0.7416	0.7448		L1-f, L2-f	0.7218	0.7353
	L0, L1-f, L2-f	0.7288	0.7430		L1-f, L1-p	0.7007	0.7176
	L0, L1-f, L1-p	0.7073	0.7253		L1-f, L2-p	0.7720	0.7862
	L0, L1-f, L2-p	0.7790	0.7939		L1-f, MIDWD	0.7428	0.7582
	L0, L1-f, MIDWD	0.7492	0.7657		L2-f, L1-p	0.7130	0.7482
3	L0, L2-f, L1-p	0.7224	0.7580	L2-f, L2-p	0.7731	0.7870	
	L0, L2-f, L2-p	0.7827	0.7977	L2-f, MIDWD	0.7449	0.7587	
	L0, L2-f, MIDWD	0.7532	0.7689	L1-p, L2-p	0.7671	0.7793	
	L0, L1-p, L2-p	0.7779	0.7912	L1-p, MIDWD	0.7287	0.7612	
	L0, L1-p, MIDWD	0.7385	0.7714	L2-p, MIDWD	0.7633	0.7883	
	L0, L2-p, MIDWD	0.7738	0.7996	L0, L1-f, L2-f, L1-p	0.7218	0.7598	
	L1-f, L2-f, L1-p	0.7147	0.7522	L0, L1-f, L2-f, L2-p	0.7823	0.8011	
	L1-f, L2-f, L2-p	0.7751	0.7935	L0, L1-f, L2-f, MIDWD	0.7530	0.7734	
	L1-f, L2-f, MIDWD	0.7463	0.7660	L0, L1-f, L1-p, L2-p	0.7767	0.7954	
	L1-f, L1-p, L2-p	0.7698	0.7878	L0, L1-f, L1-p, MIDWD	0.7393	0.7749	
5	L1-f, L1-p, MIDWD	0.7327	0.7675	L0, L1-f, L2-p, MIDWD	0.7742	0.8045	
	L1-f, L2-p, MIDWD	0.7675	0.7971	L0, L2-f, L1-p, L2-p	0.7805	0.8008	
	L2-f, L1-p, L2-p	0.7711	0.7910	L0, L2-f, L1-p, MIDWD	0.7460	0.7804	
	L2-f, L1-p, MIDWD	0.7368	0.7709	L0, L2-f, L2-p, MIDWD	0.7775	0.8046	
	L2-f, L2-p, MIDWD	0.7686	0.7945	L0, L1-p, L2-p, MIDWD	0.7723	0.8032	
	L1-p, L2-p, MIDWD	0.7623	0.7931	L1-f, L2-f, L1-p, L2-p	0.7729	0.7951	
	L0, L1-f, L2-f, L1-p, L2-p	0.7800	0.8026	L1-f, L2-f, L1-p, MIDWD	0.7388	0.7747	
	L0, L1-f, L2-f, L1-p, MIDWD	0.7457	0.7821	L1-f, L2-f, L2-p, MIDWD	0.7707	0.8004	
	L0, L1-f, L2-f, L2-p, MIDWD	0.7775	0.8078	L1-f, L1-p, L2-p, MIDWD	0.7656	0.7986	
	L0, L1-f, L1-p, L2-p, MIDWD	0.7724	0.8060	L2-f, L1-p, L2-p, MIDWD	0.7668	0.7982	
Average	L0, L2-f, L1-p, L2-p, MIDWD	0.7759	0.8076	L0, L1-f, L2-f, L1-p, L2-p, MIDWD	0.7757	0.8093	
	L1-f, L2-f, L1-p, L2-p, MIDWD	0.7688	0.8019	Average	0.7413	0.7628	

■ The scores reported are the averages from 100 randomly selected subsets of size 30 for the training set and size 10 for the test set in the 2P games, and size 45 for the training set and size 15 for the test set in the 3P games.

Size	Best Subset	Test Score
1	L2-f	0.8027
2	L0, L2-f	0.8095
3	L0, L2-f, L2-p	0.8109
4	L0, L1-f, L2-f, L2-p	0.8109
5	L0, L1-f, L2-f, L1-p, L2-p	0.8100
6	L0, L1-f, L2-f, L1-p, L2-p, MIDWD	0.8000

Table 13: Test Scores - Treatment 2P

Size	Best Subset	Test Score
1	L1-f (L2-f)	0.7867
2	L1-f, L1-p (L2-f, L2-p)	0.7903
3	L0, L2-f, L2-p	0.7927
4	L0, L1-f, L2-f, L2-p	0.7938
5	L0, L1-f, L2-f, L1-p, L2-p	0.7927
6	L0, L1-f, L2-f, L1-p, L2-p, MIDWD	0.7642

Table 15: Test Scores - 4-strategy games in 2P

Size	Best Subset	Test Score
1	L2-p	0.7716
2	L1-f, L2-p	0.7877
3	L0, L1-f, L2-p	0.7932
4	L0, L1-f, L2-f, L2-p	0.7944
5	L0, L1-f, L2-f, L1-p, L2-p	0.7926
6	L0, L1-f, L2-f, L1-p, L2-p, MIDWD	0.7901

Table 17: Test Scores - 4-strategy games in 3P

Size	Best Subset	Test Score
1	L2-p	0.7747
2	L0, L2-p	0.7787
3	L0, L1-f, L2-p	0.7815
4	L0, L1-f, L2-f, L2-p	0.7790
5	L0, L1-f, L2-f, L2-p, MIDWD	0.7765
6	L0, L1-f, L2-f, L1-p, L2-p, MIDWD	0.7728

Table 14: Test Scores - Treatment 3P

Size	Best Subset	Test Score
1	MIDWD	0.8156
2	L0, L2-f	0.8223
3	L0, L2-f, MIDWD	0.8252
4	L0, L1-f, L2-f, MIDWD (L0, L2-f, L2-p, MIDWD)	0.8237
5	L0, L1-f, L2-f, L2-p, MIDWD	0.8230
6	L0, L1-f, L2-f, L1-p, L2-p, MIDWD	0.8223

Table 16: Test Scores - 5-strategy games in 2P

Size	Best Subset	Test Score
1	L2-p	0.7778
2	L2-p, MIDWD	0.7728
3	L0, L1-f, L2-p	0.7698
4	L0, L1-f, L2-p, MIDWD	0.7642
5	L0, L1-f, L2-f, L2-p, MIDWD	0.7617
6	L0, L1-f, L2-f, L1-p, L2-p, MIDWD	0.7556

Table 18: Test Scores - 5-strategy games in 3P



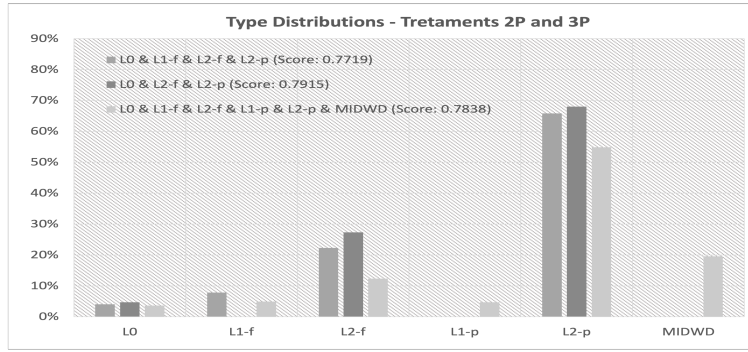


Figure 11: Type Distributions - Treatments 2P and 3P

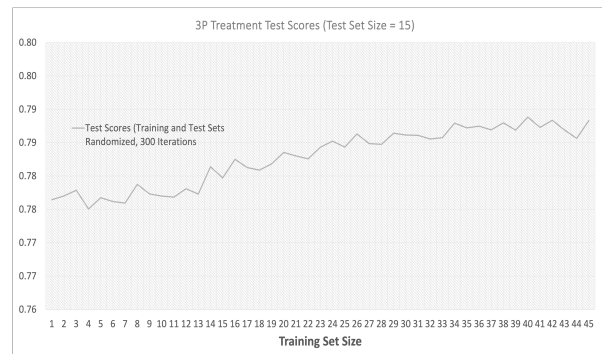
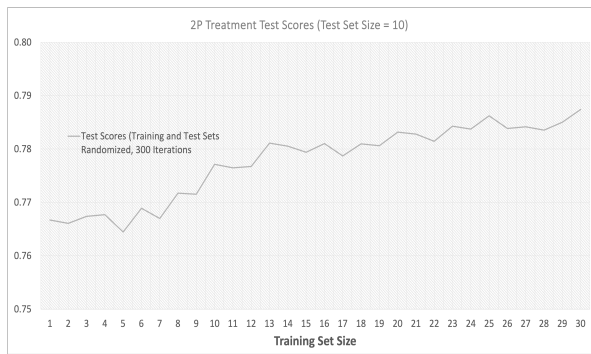
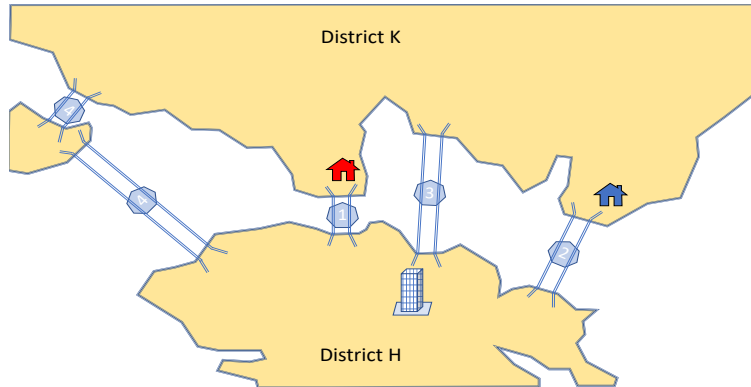


Figure 12: Test Scores with Varying Sizes of Training Set

# C Experimental Instructions - Treatment 2P

## Environment

Consider a city with two citizens, **R** and **B**. The map of the city is presented below.



The city has hired a contractor to build a bridge that connects its two districts Kowloon (K) and Hong Kong (H). A bridge is beneficial to both citizens because they both live in District K while they both work at the same office in District H. The map above indicates the location of **R**'s residence by [R's Residence], the location of **B**'s residence by [B's Residence], and the location of their office by [Office Building].

The contractor has identified four feasible locations for the bridge labelled [Bridge 1], [Bridge 2], [Bridge 3], and [Bridge 4]. The bigger the label number, the longer the bridge. The contractor's earnings depend on which bridge is built. In particular, the longer the bridge, the more the contractor earns.

Each citizen wants a short distance between his residence and the office because then he gets more points.

If the contractor builds:

- Bridge 1, then **R** will earn 200 points.
- Bridge 2, then **R** will earn 100 points.
- Bridge 3, then **R** will earn 150 points.
- Bridge 4, then **R** will earn 50 points.

If the contractor builds:

- Bridge 1, then **B** will earn 100 points.
- Bridge 2, then **B** will earn 200 points.
- Bridge 3, then **B** will earn 150 points.
- Bridge 4, then **B** will earn 50 points.

Importantly, the contractor cannot simply build any bridge he likes. In order for him to build a bridge, he needs to get permission from at least one of the citizens **R** and **B**.

## The Permission Game

You will participate in 40 rounds of decision making. In each round, you will be randomly paired with another person. One of you will be assigned to play the role of **R** and the other will be assigned to play the role of **B** (you are equally likely to be assigned the role of **R** or **B**, and the same for your partner). Your screen will tell you whether your role in the current round is **R** or **B**.

In some rounds, there will be four possible bridges, and in other rounds, there will be five possible bridges. Both the locations of possible bridges and the number of possible bridges may differ in each round. The number of points you can earn from each bridge may change in every round, both for you and the person you are paired with.

In each round, on your screen, you will see a “points table” (but not the map that corresponds to the table) that reports the number of points from each bridge for both roles, **R** and **B**. For example, the following table describes the points from the bridges located according to the map presented on the previous page:

Bridge Built	<b>R</b>	<b>B</b>
4 (the longest)	50	50
3 (2nd longest)	150	150
2 (3rd longest)	100	200
1 (the shortest)	200	100

Another example of the points table for a FIVE bridge case is as follows:

Bridge Built	<b>R</b>	<b>B</b>
5 (the longest)	100	250
4 (2nd longest)	50	50
3 (3rd longest)	250	150
2 (4th longest)	150	200
1 (the shortest)	200	100

After seeing your role and the earning table in each round, you must decide which bridge to give permission to by clicking one of the following buttons with the bridge’s number presented on your screen.

[Bridge 1]

[Bridge 2]

[Bridge 3]

[Bridge 4]

Once you have made your decision, click the “SUBMIT” button. Remember that the contractor will build the bridge with the highest number (longest length) amongst those that are chosen by you and the participant you are paired with. You will earn the points from the bridge that is built.

There is a 2-minute time limit for each round of decision making. If you do not make your choice within 2 minutes, the computer will randomly choose one of the bridges on your behalf.

After you finish making your choice in each round, the round is over, and you will not be informed about the outcome of that round until the experiment is over.

At the end of the experiment, the computer will present a table with 40 rows. Each row in the table corresponds to a round of the experiment. Each row in the table will list: (i) your choice of bridge in that round, (ii) your partner's choice of bridge in that round, (iii) the bridge that was built by the contractor in that round, and (iv) the points you earned in that round.

The computer will then randomly pick 1 round out of the 40 rounds to calculate your cash payment. So, it is in your best interest to take each round equally seriously. Your total payment in HKD will be the points you earned in that round translated into HKD via a 1:1 exchange rate plus a guaranteed HKD 40 show-up fee.

In order to get paid, you will have to fill in the receipt. The money you earn will be paid electronically via the HKUST Autopayment System to the bank account you have provided to the Student Information System (SIS). The auto-payment will be arranged by the Finance Office of HKUST and will take about 3 weeks.

Before we start the experiment for real, we will have a short Comprehension Quiz. Then you will participate in a practice round. The practice round is part of the instructions and is not relevant to your payment. The goal of the practice round is to make you familiar with the computer interface and the flow of the decisions in each round. Once the practice round is over, the computer will tell you "The official rounds begin now!"

## Comprehension Quiz

1. Suppose that the points table is given as follows: Suppose that your role is **R**. When

Bridge Built	<b>R</b>	<b>B</b>
4 (the longest)	50	50
3 (2nd longest)	150	150
2 (3rd longest)	100	200
1 (the shortest)	200	100

you give permission for Bridge 3 and **B** gives permission for Bridge 2, how many points do you earn?

- (a) 50                      (b) 100                      (c) 150                      (d) 200

2. Suppose that the points table is given as follows: Suppose that your role is **R**. When

Bridge Built	<b>R</b>	<b>B</b>
4 (the longest)	50	50
3 (2nd longest)	150	100
2 (3rd longest)	100	200
1 (the shortest)	200	150

you give permission for Bridge 1 and **B** gives permission for Bridge 3, how many points do you earn?

- (a) 50                      (b) 100                      (c) 150                      (d) 200

3. Suppose that the points table is given as follows: Suppose that your role is **B**. When

Bridge Built	<b>R</b>	<b>B</b>
5 (the longest)	150	200
4 (2nd longest)	50	100
3 (3rd longest)	250	50
2 (4th longest)	100	250
1 (the shortest)	200	150

you give permission for Bridge 2 and **R** gives permission for Bridge 5, how many points do you earn?

- (a) 50                      (b) 100                      (c) 150                      (d) 200                      (e) 250