

# Depth of Reasoning and Incentives\*

Larbi Alaoui<sup>†</sup>

Antonio Penta<sup>‡</sup>

March 5, 2013

## Abstract

We introduce a model of strategic thinking in games of initial response. Unlike standard level- $k$  models, in this framework the player's 'depth of reasoning' is endogenously determined, and it can be disentangled from his beliefs over his opponent's cognitive bound. In our approach, individuals act as if they follow a cost-benefit analysis. The depth of reasoning is a function of the player's cognitive abilities and his payoffs. The costs are exogenous and represent the game theoretical sophistication of the player; the benefit instead is related to the game payoffs. Behavior is in turn determined by the individual's depth of reasoning and his beliefs about the reasoning process of the opponent. Thus, in our framework, payoffs not only affect individual choices in the traditional sense, but they also shape the cognitive process itself. Our model delivers testable implications on players' chosen actions as incentives and opponents change. We then test the model's predictions with an experiment. We administer different treatments that vary beliefs over payoffs and opponents, as well as beliefs over opponents' beliefs. The results of this experiment, which are not accounted for by current models of reasoning in games, strongly support our theory. Our approach therefore serves as a novel, unifying framework of strategic thinking which allows predictions across games.

**Keywords:** cognitive cost – depth of reasoning – level-k reasoning – predictive game theory – strategic thinking

**JEL Codes:** C72; C91; D80.

---

\*Alaoui gratefully acknowledges financial support from the Spanish Ministry of Science and Innovation under project ECO2011-25295. We thank Ayala Arad, Ghazala Azmat, Steven Durlauf, Christian Fons-Rosen, Cristina Fuentes-Albero, Jacob Goeree, Terri Kneeland, Pablo Lopez-Aguilar, Lones Smith and especially Vincent Crawford, Nagore Iriberry and Rosemarie Nagel for their thoughtful comments.

<sup>†</sup>Universitat Pompeu Fabra and Barcelona GSE. E-mail: larbi.alaoui@upf.edu

<sup>‡</sup>University of Wisconsin-Madison. E-mail: apenta@ssc.wisc.edu

# 1 Introduction

A large body of experimental research demonstrates that individuals depart from the precepts of classical game theory in games played without clear previous strategic interaction. Studies of ‘initial responses’ suggest that individuals approach novel strategic situations by following distinct reasoning procedures of which they typically perform only a few steps. The existing literature has focused on measuring how the different procedures and depth of reasoning are distributed among subjects, but no systematic attempt has been made to understand how the depth of reasoning of individuals varies across strategic problems. Specifically, models of level- $k$  reasoning assume that individuals have an exogenous type which corresponds to the number of rounds of iterated reasoning they perform. In these models, a level-0 individual represents a non-strategic type that follows some exogenously specified behavior, while a level-1 individual best responds to level-0, and so forth.<sup>1</sup> By taking individuals’ depth of reasoning as exogenous, these models are silent over how the players’ cognitive process and actions vary with the strategic setting. From an empirical viewpoint, it may be that players change their level of play as their incentives to reason and beliefs over their opponents’ cognitive constraints vary. Verifying these conjectures and developing an explicit model of the reasoning process that delivers predictions across games would significantly increase the power of this approach.

In this paper, we introduce a framework in which players’ depth of reasoning is endogenously determined as resulting from a procedure that relates individuals’ cognitive abilities to the payoffs of the game. Behavior in turn follows from the individual’s depth of reasoning and his beliefs about the reasoning process of the opponent. Thus, in our approach, payoffs not only affect individual choices in the traditional sense, but they also shape the cognitive process itself. We then present an experimental test of our theory. The experimental results reveal that individuals change their behavior in a systematic way as payoffs and opponents change, thereby confirming that incentives and beliefs play an important role in determining the agents’ depth of reasoning and level of play. Moreover, these findings are consistent with our theoretical predictions and strongly support the model.

The fundamental feature of our framework is that players act as if they weigh the incremental value of additional rounds of reasoning against an incremental cost of learning more about the game from introspection. While the cognitive cost is exogenous, the ‘value of reasoning’ is connected to the game payoffs. In this model, increasing the stakes of the game provides individuals with stronger incentives to reason. These increased incentives may induce them to perform more rounds of reasoning. But depth of reasoning need not coincide with the sophistication of the chosen action. When facing opponents that they perceive to be more sophisticated

---

<sup>1</sup>Level- $k$  models were first introduced by Nagel (1995) and Stahl and Wilson (1994, 1995). Camerer, Ho and Chong (2004) propose the closely related ‘cognitive hierarchy’ model, in which level- $k$  types respond to a distribution of lower types. Level- $k$  models have been extended to study communication (Crawford, 2003), incomplete information (Crawford and Iriberry, 2007) and other games. Crawford, Costa-Gomes and Iriberry (2012) provide a thorough survey of this literature and Fudenberg (2010) discusses the importance of this approach to the development of predictive game theory. For recent theoretical work inspired by these ideas, see Strzalecki (2010), Kets (2012) and Kneeland (2013).

than themselves, subjects play according to their own cognitive bound. But when facing less sophisticated opponents, then they play according to less rounds of introspection than their actual cognitive bound. We note that the notion of playing a more sophisticated opponent is natural in this setting, thereby resolving a well-known conceptual difficulty of the level- $k$  approach. Our model further predicts that individuals follow (weakly) more sophisticated behavior when the opponents' incentives to reason are increased, unless their own cognitive bound is binding. Depending on the game, these predictions on the depth of reasoning and behavior translate to stochastic dominance relations in the distributions of actions as incentives, payoffs and beliefs over the opponents are varied. We describe the predicted stochastic dominance relations in detail in Section 3.

A cost-benefit approach to modeling the reasoning process has several advantages. From a methodological viewpoint, it bridges the study of strategic thinking with standard economic concepts. It also provides a tractable mechanism for understanding the conceptually complex interaction between a player's reasoning about the game, his reasoning about the opponent's reasoning procedure and his choice of action. Furthermore, the rich set of testable predictions that it delivers do not rely on assumptions on the specific shapes of the cost and value of reasoning functions. Since strategic reasoning is not a conventional domain of analysis, obtaining these results with minimal imposed structure is an important feature of our approach.<sup>2</sup> Moreover, in our model, the level of reasoning according to which an agent plays is the endogenous outcome of a reasoning process, and not a fixed parameter. By making explicit an appealing feature of level- $k$  models that play follows from a reasoning procedure, our framework serves to attain a deeper understanding of the underlying mechanisms of this approach. Our framework further enriches level- $k$  theory by disentangling the rounds of introspection that a player performs from the rounds he believes his opponents perform.

Our experiment tests the predictions that hold with full generality of the model. It also serves the broader purpose of testing the conjecture that players vary the number of rounds of reasoning they perform as their incentives and beliefs over opponents change.<sup>3</sup> Our design consists of varying the players' beliefs and their incentives. We consider two different ways of changing the agents' beliefs over their opponents cognitive abilities. In both cases, we divide the subjects into two groups whose labels are perceived to be informative about game theoretic sophistication. In the first case, we separate the subjects into two groups by degrees of study. In the second, subjects are required to take a test of our design, and are then separated by

---

<sup>2</sup>Whereas in the present paper we introduce our framework and test its core implications, in Alaoui and Penta (2013) we pursue an axiomatic approach to players' reasoning, in which the cost-benefit analysis emerges as a representation. Besides uncovering the fundamental underpinnings of the approach, the axioms enable us to impose structure on the functional forms. We further show that, endowed with this added structure, our model is highly consistent with all the well-known static treasures of Goeree and Holt (2001).

<sup>3</sup>With respect to the importance of beliefs on players' actions, a recent experiment by Agranov, Potamites, Schotter and Tergiman (2012) makes the simple but important point that beliefs do change the average number of rounds performed, in a standard beauty contest. Palacios-Huerta and Volij (2009) explore a related point in the dynamic context of the centipede game. They analyze how the number of rounds of backward induction performed by subjects vary when the level of sophistication of the opponent is changed.

their score, which can either be ‘high’ or ‘low’. We then use these labels to vary agents’ beliefs over their opponents’ cognitive constraints. These changes serve to test the model’s predictions that agents play according to a lower depth of reasoning when playing against opponents they take to be less sophisticated. Our theory also allows players to not only take into account the (perceived) sophistication of the opponent, but also the opponent’s belief over the player’s own sophistication. To account for these higher order beliefs effects, we administer treatments in which subjects classified under a label play against the action that subjects from the other label have played against each other.

To test whether players respond to increased incentives of doing more rounds of introspection in the manner that is predicted by our model, we increase their reward for being ‘correct’ in their reasoning. We also test the prediction that an increase in opponents’ incentives in itself leads subjects to perform more rounds. In particular, we administer treatments in which subjects playing the high reward game play against an action chosen by opponents playing the low reward game against each other. We then compare the distributions of the chosen actions across these different treatments. Our results are consistent with the predictions that subjects play according to more rounds of introspection when stakes are increased and when opponents are believed to be more sophisticated. The results are also in line with the predictions over higher order beliefs effects. Our model further allows for the analysis of the experiment’s more complex observed patterns. In particular, the observed shifts in distributions as beliefs over opponents are changed are more pronounced when players’ incentives to reason are weaker. These findings are indicative of an interaction between changes in incentives and changes in beliefs over opponents that is within the scope of our model. Overall, the experimental results of the main treatments strongly support our theoretical predictions.

Lastly, we provide additional treatments to explore more subtle implications of our theory. These designs are more complex and cognitively demanding for the subjects, and some rely on our theory more than they directly test it. For instance, in order to identify individuals’ beliefs about the cognitive abilities of their opponents, our model suggests that individuals’ cognitive costs would have to be sufficiently low. While costs of reasoning are in principle exogenous, we lower them by exposing subjects to a game theoretic ‘tutorial’ of our game, after having administered the main treatments. The subjects then play the game against other subjects who were also given the tutorial, as well as against actions that were chosen in some of the baseline (pre-tutorial) treatments. In further support of our framework, the effects of the game theory tutorial are comparable to those observed in the baseline treatments, when the incentives to reason are increased.

The experiment therefore shows that individuals change their behavior in a systematic manner that is not endogenized by existing models of strategic reasoning, but that is strongly consistent with our theoretical predictions. By jointly accommodating this set of results in a tractable and intuitive way, our cost-benefit approach to modeling the reasoning process is both empirically relevant and instrumental to a better understanding of reasoning in games. The scope and predictive power of this approach do not depend on parametric restrictions or

added structure to the core model. Our model therefore serves as a unifying framework for the analysis of behavior in games of initial response.

The paper is structured as follows. For clarity of exposition, we introduce the experimental design in Section 2 and the theoretical model and its predictions in Section 3. Section 4 analyzes the empirical results, Section 5 focuses on additional treatments, and Section 6 concludes.

## 2 Experimental Design

We first describe the experiment before presenting the theoretical model. In essence, the experiment is designed not only to test whether individuals play differently when their incentives and beliefs about opponents change, but also to analyze the direction in which their actions change. Moreover, we aim to disentangle whether their action is dictated by their cognitive constraints, given their incentives, or by their beliefs over their opponents' cognitive constraints. The baseline game remains the modified 11-20 game throughout:

The subjects are matched in pairs. Each subject enters an (integer) number between 11 and 20, and always receives that amount in tokens. If he chooses *exactly one less* than his opponent, then he receives an extra 20 tokens. If they both choose the same number, then they both receive an extra 10 tokens.

This game is a variation of Arad and Rubinstein's (2012) '11-20' game, the distinction being that the original version does not include the extra reward in case of a tie (where one token was worth five euro cents in our experiment). We note that aside from our use of a similar game, our objectives and experimental design are radically different.

As argued by Arad and Rubinstein, the 11-20 game presents a number of advantages in the study of level- $k$  reasoning, which are inherited by our modified version. We recall here the most relevant to our purposes. First, using level- $k$  reasoning is natural, as there are no other obvious focal ways of approaching the game. The competing alternative of guessing the unique pure-strategy equilibrium seems far from self-evident, and would be difficult to see without going through at least a few steps of iterated reasoning. Secondly, the level-0 specification is intuitively appealing and unambiguous, since choosing 20 is a natural anchor for an iterative reasoning process. Moreover, it is the unique best choice for a player who ignores all strategic considerations. Thirdly, there is robustness to level-0 specification, in that the choice of 19 would be the level-1 strategy for a wide range of level-0s, including the uniform distribution over the possible actions. Lastly, best-responding to any level- $k$  is simple: level-1 plays 19, level-2 best responds to 19 by playing 18, and so on. Since we do not aim to capture cognitive limitations due to computational complexity, having a simple set of best responses is preferable in this case.

In addition to these points, our modification of the game leads to another useful feature for our objectives. By introducing the extra reward in case of tie, the best response to 11 is 11,

and not 20, as in the version of Arad and Rubinstein. Thus, our modification breaks the cycle in the chain of best responses, which enables us to assign one specific level of reasoning to each possible announcement (with the exception of 11, which corresponds to any level equal to 9 or higher). Action 19 can only be a level-1 strategy, 18 can only be a level-2 strategy, and so forth for every  $k$  up to  $k = 8$ . In the original 11-20 game, action 19 could have been played by a level 1, but also by a level-11, level-21, or other ‘high’ levels (levels of form  $10n + 1$ ). Although levels-11 and above appear to be uncommon, it is crucial that these cycles be avoided here. One of the main hypotheses that we aim to test is whether an increase in players’ incentives would shift the distribution of level- $k$ ’s toward higher  $k$ ’s, but this hypothesis could not be falsified in the presence of such cycles.

The subjects of the experiment were 120 undergraduate students from different departments at the Universitat Pompeu Fabra (UPF), in Barcelona. Each subject played twice every treatment described in Sections 2.1 and 2.2, and summarized in Table 1. They also played a subset of the additional treatments described in Section 5. These treatments are all based on the modified 11-20 game. We provide the exact sequences of treatments used in Appendix A.2.

Each subject was anonymously paired with a new opponent after every iteration of the game. To focus on initial responses and to avoid learning from taking place, the subjects did not observe their payoffs after their play. They only observed their earnings at the end of the session. Moreover, subjects were paid randomly, and therefore did not have any mechanism for hedging against risk by changing their actions.<sup>4</sup> Specifically, they were paid once for each set of six iterations, and the randomization occurred inside this set (once at random for the first six iterations, once for the next six iterations, and so forth). As an additional control for order effects, the order of treatments was randomized. Furthermore, since subjects played the same treatments twice during a session, we can compare play for each treatment through equality of distribution tests. Lastly, subjects received no information concerning their opponents’ results. This serves to avoid that subjects focus on goals other than monetary incentives, such as defeating the opponent or winning for its own sake. The instructions of the experiment were given in Spanish; the English translation and the details on the pool of subjects, the earnings and the logistics of the experiment are in Appendix A.

## 2.1 Changing beliefs about the opponents

We consider two different classifications of subjects, an *exogenous classification* and an *endogenous classification*, each with 3 sessions of 20 subjects. In the exogenous classification, subjects are distinguished by their degree of study. Specifically, in each session of the experiment, 10 students are drawn from the field of humanities (humanities, human resources, and translation), and 10 from math and sciences (math, computer science, electrical engineering, biology

---

<sup>4</sup>These devices are standard in the literature that focuses on ‘initial responses’, where the classical equilibrium approach is hard to justify. See, for instance, Stahl and Wilson (1994, 1995), Costa-Gomes, Crawford and Broseta (2001) and Costa-Gomes and Crawford (2006). For an experimental study of equilibrium in a related game, see Capra, Goeree, Gomez and Holt (1999).

Treatment	Own label	Opponent's label	Own payoffs	Opponent's payoffs	Replacement of opponent's opponent
Homogeneous [A]	$I$ ( $II$ )	$I$ ( $II$ )	Low	Low	No
Heterogeneous [B]	$I$ ( $II$ )	$II$ ( $I$ )	Low	Low	No
Replacement [C]	$I$ ( $II$ )	$II$ ( $I$ )	Low	Low	Yes
Homogeneous-high [A+]	$I$ ( $II$ )	$I$ ( $II$ )	High	High	No
Heterogeneous-high [B+]	$I$ ( $II$ )	$II$ ( $I$ )	High	High	No
Replacement-high [C+]	$I$ ( $II$ )	$II$ ( $I$ )	High	High	Yes

Table 1: Treatment summary: Label  $I$  refers to ‘math and sciences’ or to ‘high’ subjects, and label  $II$  refers to ‘humanities’ or to ‘low’ subjects.

and economics). They are aware of their own classification when beginning the experiment, and are labeled as ‘humanities’ or ‘math and sciences’. In the endogenous classification, there is no restriction on the pool of subjects. Moreover, the subjects are not informed about the field of study of the other players. Before playing the game, however, they are required to take a test of our design. Based on their performance on this test, each student is either labeled as ‘high’ or ‘low’, and is shown his own label before playing the game. We defer the description of this test to Section 2.3.1 (see also Appendix B).

These classifications allow us to change subjects’ beliefs about their opponents. In each treatment, the subjects are given information concerning their opponents. They play the baseline game against someone from their own label (homogeneous treatment [A]) and against someone from the other label (heterogeneous treatment [B]). For instance, for the exogenous classification, a student from math and sciences (resp., humanities), is told in homogeneous treatment [A] that his opponent is a student from math and sciences (humanities) as well. In heterogeneous treatment [B], he is told that the opponent is a student from humanities (math and sciences). Identical instructions are used for the endogenous classification, but with ‘high’ and ‘low’ instead of ‘math and sciences’ and ‘humanities’, respectively.

Suppose, for the sake of illustration, that there are two cognitive types of subjects, consisting of those with higher game theoretical sophistication and those with lower game theoretical sophistication. For our purposes, it suffices that players *believe* that there is a meaningful difference between these two types, and that the labels we use in the two classifications are informative about their opponents’ type (see Section 2.3). We hypothesize that subjects in the exogenous classification associate the ‘maths and sciences’ label and ‘humanities’ label with higher and lower game theoretical sophistication, respectively, and that subjects in the endogenous classification associate the ‘high’ and ‘low’ labels with higher and lower game theoretical sophistication, respectively.<sup>5</sup> Then, when playing homogeneous treatment [A] compared to heterogeneous treatment [B], subjects change from believing they are playing against one type of player to another.

Treatments [A] and [B] are designed to test whether the behavior of the subjects varies with the sophistication of the opponent. The next treatment is designed to test whether the subjects believe that (or are aware of) the behavior of their opponents also changes when they

<sup>5</sup>This assumption is not required, as it is revealed by the data. Since our results are consistent with this interpretation, we maintain it in this discussion.

face opponents of different levels of sophistication. To do so, we consider replacement treatment [C]. A ‘math and sciences’ subject, for instance, is given the following instructions: “[...] two students from humanities play against each other. You play against the number that one of them has picked.” The reasons for using this exact wording are discussed in Section 2.3.2.

## 2.2 Changing incentives

We next consider a second dimension that would entail a change in players’ chosen actions, according to our framework. In particular, we aim to test the central premise of our theoretical model, that players may perform more rounds of introspection if they are given more incentives to do so. To do this, we change the baseline game in the following way: rather than winning an extra 20 tokens for choosing the number precisely one below the opponent, subjects win an extra 80 tokens. The rest of the game payoffs remain the same. It is immediate that this change does not affect the level- $k$  actions, irrespective of whether the level-0 is specified as 20 or as the uniform distribution. It only increases the rewards for players who stop at the ‘correct’ round of reasoning.

We consider three treatments for this ‘high payoff game’: homogeneous treatment [A+], heterogeneous treatment [B+], and replacement treatment [C+]. These treatments are equivalent to treatments [A], [B] and [C], respectively, but with the increased reward for undercutting of 80 tokens. We then compare an agent’s play under lower payoffs to his play under higher payoffs, by comparing [A] to [A+], [B] to [B+] and [C] to [C+]. We also compare treatments [A+], [B+] and [C+] in an analogous way to the comparison between treatments [A], [B] and [C].

This concludes our discussion of the main treatments. The theoretical framework introduced in the next section makes predictions on the change in the distribution of actions chosen by players of both groups, for each of the six treatments in Table 1. These predictions are summarized in Table 2 of Section 3. We later consider additional treatments, including treatments analogous to replacement treatment [C], but where the payoffs are replaced: a subject with a high payoff plays against the number chosen in a low payoff treatment. A description of these treatments is presented in Section 5.

## 2.3 Experimental Design: Discussion

### 2.3.1 Designing the Group Classification: Demarcation and Focality.

In order to vary subjects’ beliefs about the opponents, we divide the pool of subjects into two labeled groups. We then change subjects’ beliefs about the opponents by changing the opponent’s group in the different treatments. This approach requires two features. The first is *demarcation*: the types in the classifications must be perceived as being sufficiently distinct from each other in terms of their strategic sophistication. The second is *focality*: since subjects’ behavior depends not only on their beliefs but also on their beliefs about their opponents’



beliefs, it is important that the two types share sufficient agreement about the way the two types differ. That is, they should ‘commonly agree’ over their relative sophistication. The two classifications we consider have been chosen to guarantee that these properties hold.

**The Exogenous Classification.** The exogenous classification exploits the intuitive, albeit vague, view that ‘math and sciences’ students are regarded as more accustomed to numerical reasoning than ‘humanities’ students. Furthermore, the specific degrees of study used to populate the ‘math and sciences’ group are commonly viewed as being the most selective degrees at UPF, and require the highest entry marks.<sup>6</sup> We would therefore expect the subjects to believe the ‘math and sciences’ group to be comparatively more sophisticated in game theoretical reasoning than the ‘humanities’ group. However, the subjects are not primed into shaping specific beliefs about either particular group.

**The Endogenous Classification and the Test.** Since the exogenous classification is based on labels that are salient for the subjects in the pool, it can be expected to ensure both demarcation and focality. But it does not allow us to fully control the agents’ beliefs about these labels. For this reason, we also introduce the endogenous classification, where students are classified solely based on their performance in a test of our design. The goal of the test is thus twofold. It sorts subjects into two groups, and, by labeling the scores obtained by subjects as ‘high’ or ‘low’, the test itself forms the agents’ beliefs over the content of these labels.

The main objective of the test is to convince subjects that the result of the test is informative about their opponents’ game theoretical sophistication. To do so, we ensure that our test questions appear difficult to solve, and that subjects would be likely to infer that an individual of higher sophistication would respond better to the questions.<sup>7</sup>

The cognitive test takes roughly thirty minutes to complete, and consists of three questions. These are all single-agent questions, in that they are not pitted against each other. Rather, subjects are asked to provide the correct answer. Scores are assigned to the answers to each of these tasks, and a formula then takes a weighted average. Those who have a score above the median are labeled ‘high’, and the others are labeled ‘low’. Subjects do not see their numerical grade, but they are told whether they are labeled ‘high’ or ‘low’. (Details of the test are contained in Appendix B).

The three questions are as follows. In the first question, subjects have nine attempts to play a variation of the board game *mastermind*, in which the aim is to deduce a hidden pattern through a sequence of guesses. This game requires skill at logical inference to complete successfully. In the second question, they are given a typical *centipede game* of seven rounds.

---

<sup>6</sup>These views emerged from informal conversations with students. They are confirmed by the admission scores, used to select the students admitted in the various fields. These scores can be found at: <http://www.elpais.com/especial/universidades/titulaciones/universidad/universidad-pompeu-fabra/45/nota-corte/>.

<sup>7</sup>Another possible consequence of being assigned a ‘high’ or a ‘low’ label may impact an individual’s self-esteem, which may in turn impact his subsequent performance in playing the modified 11-20 game. This concern, however, is tangential to the aim of our experiment. Our objective is not to identify the number of rounds of introspection a subject performs in a single game, but how his actions vary across treatments.

In the third game, the agents are given a lesser known *pirates game*. As with the centipede game, the pirates game can be solved by backward induction. The solution is more difficult to attain, however, in that it involves additional computation in determining how players' strategy profiles map into outcomes. For the latter two games, subjects are not asked how they would play; rather, they are asked how "infinitely sophisticated and rational agents, who each want to get as much money as possible" would play.

These three questions are challenging to most, and performing well is arguably informative of a subject's sophistication. More importantly, for our purposes, it seems fair to assume that the *subjects* believe the score to be a strong indicator of game theoretical sophistication. Arguably, this would be the case for subjects who recognize that the 11-20 game has a recursive pattern reminiscent of the structure of the problems in the test, which appears plausible for players who would perform at least one round of introspection. Furthermore, while the last two games are close enough to the 11-20 game to require a similar kind of analysis, they are sufficiently different that the feedback provided in the form of the 'high' or 'low' score should limit the learning about how to play the 11-20. It is precisely to avoid information leaking from the test to the proper experiment that we have not included the beauty contest, the traveler's dilemma or other closely related games.

### 2.3.2 Testing for Effects of Higher Order Beliefs

Suppose that label *I* denotes the group perceived by the subjects as relatively more sophisticated, and that label *II* denotes the group perceived as relatively less sophisticated. We would then expect label *I* subjects to play (weakly) lower numbers in treatment [A] than in treatment [B], and the label *II* subjects to play (weakly) lower numbers in treatment [B] than [A]. Whereas these treatments test whether subjects' beliefs affect their choices in the game, it is difficult to establish from these treatments alone whether the subjects are aware that their opponents' beliefs may affect their choices. For instance, subjects may expect, as we do, that a label *II* subject plays differently against another label *II* than against label *I*. This would indicate that they have an understanding that the label *II* subject may not necessarily play according to his cognitive limit, and that his play may vary.

The objective of treatment [C] is precisely to test for the degree to which the subjects have a well-formed model of their opponents' reasoning, and whether it is consistent with our theoretical predictions. The precise wording of treatment [C] is designed to pin down the entire hierarchy of beliefs, since any potential ambiguity may lead to a misinterpretation of the results. For instance, the full description that a math and sciences student is given concerning his opponent in treatment [C] is: "[...] two students from humanities play against each other. You play against the number that one of them has picked." It is therefore clear that he is playing a humanities playing a humanities subject, who himself is playing a humanities subject, and so forth. Any ambiguity that allowed players to believe that one of them believes (at some high level) that one of them is a student from math and sciences could result in a sophisticated subject behaving as a less sophisticated one, invalidating the identification of the types.

### 2.3.3 Choice of the Baseline Game

While the modified 11-20 game is particularly suitable to our purposes, our model applies to a wide spectrum of games. We discuss here some related games to which our framework appears especially relevant.

Basu’s (1994) traveler’s dilemma is perhaps the closest to ours. In this two-player game, each agent reports a number between 2 and 100, and both players receive the lowest of the two numbers chosen. In addition, if the players report different numbers, then the one who reports a higher number pays a penalty of 2, and the one with the lowest receives a reward of 2. This game shares appealing features of the modified 11-20 game. The main difference is that it is sufficient in the traveler’s dilemma to undercut the opponent to receive the additional reward, rather than to choose exactly the right action. The modified 11-20 game therefore leads to a more precise mapping between the agent’s action and his beliefs. Moreover, if agents had social preferences such as altruism or fairness, it would not be an issue in the modified 11-20 game. This is because, independent of the agents’ preferences over the final outcomes, the optimal choice would still require players to identify their opponents’ action.<sup>8</sup>

A central position in the literature on level- $k$  is occupied by Nagel’s (1995) *beauty contest* game, also known as the *p-guessing game*.<sup>9</sup> In this game,  $n$  players are asked to report a number between 0 and 100, and the player whose number is closest to a fraction  $p \in (0, 1)$  of the average report wins a monetary prize. A remarkable finding of the literature is that players’ reports exhibit a regular pattern concentrated around specific numbers, and that the position of these spikes shifts down as  $p$  decreases. This evidence has been interpreted as suggesting that players approach games of this kind by thinking in steps, which has been a key motivation to the development of level- $k$  theories. Importantly, the large strategy space of this game allows testing whether players reason according to some form of iterated reasoning.<sup>10</sup> Our goal, however, is different: here we maintain that players follow an iterated reasoning procedure, and we ask whether their depth of reasoning varies with their beliefs about the sophistication of the opponents, and with incentives to reason more deeply about the game. A number of features of the beauty contest makes it less suitable than the modified 11-20 game for this task. For instance, level- $k$  reasoning is not necessarily the only ‘focal’ form of reasoning, and, for related reasons, the beauty contest does not present an obvious specification of the level-0 action.<sup>11</sup> This last point is particularly problematic because the level-1 action is highly sensitive to the specification of level-0 in this game. Furthermore, a natural setting for our experiment is one with few players, as it allows for greater ease in changing beliefs over

---

<sup>8</sup>A traveler’s dilemma in which payoffs are varied can be found in Goeree and Holt (2001). The results of this game fit well with our theoretical predictions; see Alaoui and Penta (2013) for details.

<sup>9</sup>Variations of this game have been studied, among others, by Ho, Camerer and Weigelt (1998) and Bosch-Domènech, García-Montalvo, Nagel and Satorra (2002).

<sup>10</sup>For studies that focus more directly on the cognitive process itself, see Agranov, Caplin and Tergiman (2012), and the recent works by Bhatt and Camerer (2005), Coricelli and Nagel (2009), and Bhatt, Lohrenz, Camerer and Montague (2010), which use fMRI methods and find further support for level- $k$  models.

<sup>11</sup>For a thorough description of these different thought processes, see Bosch-Domènech, García-Montalvo, Nagel and Satorra (2002).

opponents, and opponents’ opponents. But with a small number of players, the best response function in the  $p$ -guessing game is difficult to compute. A player must take into account the impact of his own number on the average, which is not as immediate a calculation as computing the best response in our game (see Grosskopf and Nagel 2008). Notwithstanding these factors, we expect the comparative statics predictions of our theoretical model to hold for the beauty contest as well.

In another important game, introduced by Costa-Gomes and Crawford (2006, CGC), the players’ objective is to guess a target that depends on their opponent’s guess multiplied by a constant. Because of the large strategy space, this game inherits the appealing properties of the beauty contest. Like our game, however, it is a two-player game with simple best response functions. It further allows a separation between different types of reasoning processes, such as level- $k$  reasoning and iterated dominance, which is not allowed by the standard beauty contest. More importantly, subjects in CGC play a sequence of games in which both the strategy space and the target are varied. Such sequences of responses yield strategic ‘fingerprints’ that allow CGC to identify individuals’ types of reasoning. As discussed earlier, our objective is distinct, as we do not aim to separate level- $k$  from other forms of reasoning.

The modified 11-20 game is apt for focusing on the novel implications of our theoretical model, while minimizing the impact of confounding factors. But the games described here arguably fall within the scope of our model, and can therefore serve to better identify the cognitive procedure associated with observed behavior. Moreover, our model can be instrumental to the analysis of these games, particularly as players’ incentives and beliefs are varied. As discussed in Alaoui and Penta (2013), our theoretical predictions are highly consistent with observed behavior in a number of experiments.<sup>12</sup> In addition to applying our experimental design to these alternative games, interesting directions for future research include enriching the experiment along the lines of CGC, so as to elicit profiles of responses to better identify individual cognitive processes.

### 3 A Theory of Endogenous Level- $k$ Reasoning

This section introduces our model of endogenous level- $k$  reasoning. We take as given that players approach the game with an iterated reasoning process, and endogenize their depth of reasoning. The number of steps they take is a function of their game theoretical sophistication and the payoff structure of the game. We also endogenize the players’ chosen actions, which depend both on the number of steps of reasoning they perform and on their beliefs about the cognitive abilities of their opponent. We then show that this theory delivers testable predictions for our experiment. In particular, the model implies first-order stochastic dominance relations between the distributions of the actions played by individuals of a given ‘label’ in the different treatments of the experiment.

---

<sup>12</sup>For instance, the predictions of our theory in the five static games considered in Goeree and Holt (2001) are consistent with their experimental findings.

### 3.1 Steps of Reasoning

Consider a finite two-players game with complete information,  $G = (A_i, u_i)_{i=1,2}$ , where  $A_i$  is the (finite) set of actions of player  $i$ ,  $A_1 = A_2$ , and  $u_i : A_1 \times A_2 \rightarrow \mathbb{R}$  is player  $i$ 's payoff function. Let  $G$  be such that the (pure strategy) best response correspondence  $BR_i : A_j \rightarrow A_i$ , defined as

$$BR_i(a_j) = \arg \max_{a_i \in A_i} u_i(a_i, a_j) \text{ for each } a_j \in A_j,$$

is single-valued. Exploiting the symmetry of the action space and the single-valuedness of  $BR_i$ , we assume that players' reasoning about the game is represented by sequences  $\{a_1^k\}_{k \in \mathbb{N}}$ ,  $\{a_2^k\}_{k \in \mathbb{N}}$  such that  $a_1^0 = a_2^0$ , and  $a_i^{k+1} = BR_i(a_j^k)$  for each  $k \in \mathbb{N}$ . We view  $a_i^0$  to be the action that player  $i$  would play by default, without any strategic understanding of the game. As player  $i$  performs the first step of reasoning, however, he becomes aware that his opponent could play  $a_j^0$ , and thus considers playing  $a_i^1 = BR_i(a_j^0)$ . Similarly, as player  $i$  advances from step  $k-1$  to step  $k$ , he realizes that his opponent may play  $a_j^{k-1}$ , in which case the best response would be  $a_i^k = BR_i(a_j^{k-1})$ .

We interpret the steps of reasoning as 'rounds of introspection'. In our model, players are not boundedly rational in the sense of failing to compute best responses. Rather, players are limited in their ability to conceive that the opponent may perform the same steps of reasoning. As an illustration, consider the modified 11-20 game used in the experiment. For any player  $i$ , action  $a_i^0 = 20$  is a natural action for a level-0 player, as it is the number that a non-sophisticated player would report if he ignored all strategic considerations.<sup>13</sup> If player  $i$  exerts some cognitive effort and performs the first step of the reasoning process, then he realizes that the opponent may approach the game in the same way, in which case the best response is  $a_i^1 = 19$ . Similarly, if player  $i$  performs one more round of introspection, he realizes that his opponent might also perform the same reasoning, and play 19. In that case, the best response is  $a_i^2 = 18$ . This reasoning continues until he reaches 11, in which case the best response remains at 11.<sup>14</sup> However, this does not necessarily mean that player  $i$  would play according to the number of steps of reasoning he has performed. His actions depends on his beliefs about player  $j$ 's cognitive abilities. For instance, if player  $i$  has performed three steps of reasoning then he does not play 17 if he believes that  $j$  has not performed two steps of reasoning.

This interpretation of level- $k$  is related to that of Crawford and Iriberri (2007) and other models of level- $k$  reasoning. In these frameworks, a player's  $k$  identifies both the action that he plays and his 'understanding of the game'. Here, we distinguish the two. The action chosen by a player is a function of both the highest  $k$  he has reached and his beliefs about the number of rounds performed by his opponents. Denoting by  $\hat{k}_i$  player  $i$ 's *cognitive bound*, our model endogenizes both the cognitive bound  $\hat{k}_i$  and the  $k_i$  according to which he plays. The former is only a function of player  $i$ 's incentives and cognitive abilities, while the latter is also a function

<sup>13</sup>As discussed in Section 2, different specifications of the level-0 (including the uniform distribution) would not affect the analysis. For simplicity, we only consider  $a_i^0 = 20$  here.

<sup>14</sup>Formally,  $BR_i(a_j)$  is single-valued for any pure strategy  $a_j$ , and is equal to  $BR_i(a_j) = \max\{a_j - 1, 11\}$  for  $a_j \in \{11, \dots, 20\}$ . The action associated with any given  $k \geq 0$  is  $a_i^k = \max\{20 - k, 11\}$ .

of his beliefs about the cognitive abilities of the opponent.<sup>15</sup> Player  $i$  may have cognitive bound  $\hat{k}_i$  but play action  $a_i^{k_i}$ , where  $k_i \leq \hat{k}_i$ , if he thinks that  $a_i^{k_i}$  is preferable to  $a_i^{\hat{k}_i}$ .

### 3.2 Cognitive Costs and Value of Introspection.

The model we propose for endogenizing the steps of reasoning taken by players is based on a cost-benefit analysis. Performing additional rounds of reasoning entails incurring a cognitive cost. While these costs reflect a player’s cognitive ability, which we view as exogenous, we assume that the benefits of performing an extra step of reasoning depend on the payoff structure of the game.

We stress that we do not view this ‘cost-benefit’ analysis as an optimization problem *actually* solved by the agent, but rather as a modeling device to represent a player’s reasoning about the game. We hypothesize that agents’ understanding of the game varies systematically with the payoff structure. To the extent that players’ understanding of the game exhibits this form of consistency, it can be modeled *as if* the cognitive bound  $\hat{k}_i$  results from a cost-benefit analysis. In Alaoui and Penta (2013) we provide an axiomatic foundation to this approach, deriving the cost-benefit representation from primitive assumptions on the player’s reasoning process, explicitly modeled as a Turing machine.

#### 3.2.1 Individual Understanding of the Game.

Formally, we assume that the value of doing extra steps of reasoning only depends on the payoff structure of the game. Fixing the game payoffs, we define function  $v_i : \mathbb{N} \rightarrow \mathbb{R}_+$ , where  $v_i(k)$  represents  $i$ ’s value of doing the  $k$ -th round of reasoning, given the previous  $k - 1$  rounds. The cognitive ability of agent  $i$  is represented by a cost function  $c_i : \mathbb{N} \rightarrow \mathbb{R}_+$ , where  $c_i(k)$  denotes  $i$ ’s incremental cost of performing the  $k$ -th round of reasoning. The following assumptions on  $c_i$  and  $v_i$  are maintained throughout.

**Condition 1** *Maintained assumptions on the Cost and Value of Reasoning:*

1. **Cost of Reasoning:**  $c_i(0) = 0$  and  $c_i(k) \geq 0$  for every  $k \in \mathbb{N}$ .
2. **Value of Reasoning:**  $v_i(k) \geq 0$  for every  $k \in \mathbb{N}$ .

Condition 1 entails minimal restrictions on the cost and value of reasoning functions. In particular, it contains no assumptions about their shape (their monotonicity, convexity, etc.). Maintaining this level of generality allows us to focus on the essential features of our approach and to capture different kinds of plausible cost functions. For instance, in the modified 11-20

---

<sup>15</sup>In the following we introduce different notions of level  $k$ . Without a superscript,  $k_i$  refers to player  $i$ ’s actual  $k$ . When both a subscript and a superscript are present, the symbol denotes a belief: the subscript indicates the player that a particular  $k$  refers to, and the superscript indicates the player holding the belief. For instance,  $k_j^i$  denotes  $i$ ’s belief about  $j$ ’s behavioral  $k$ . Moreover, the  $k$  notation refers to behavior while  $\hat{k}$  notation refers to cognitive bounds (e.g.  $\hat{k}_i$  is  $i$ ’s cognitive bound, and  $\hat{k}_j^i$  is  $i$ ’s belief about  $j$ ’s cognitive bound.)

game, a player who understands the inductive structure of the problem would have a non-monotonic cost  $c_i$ . His first rounds of reasoning would be cognitively costly but those following the understanding of the recursive structure would not be, as described in Example 1.

Whereas the cost function represents the cognitive ability of the player, function  $v_i$  represents the value of performing each extra round of reasoning in the game. The assumption that  $v_i(k) \geq 0$  captures the idea that, net of its cost, a deeper understanding of the game is never detrimental to the agent, who can at worst ignore the extra insight of each further step of reasoning. We take the value of reasoning to be purely instrumental to informing the player's choice of action. The  $v_i$  function thus depends only on the payoff structure of the game, with the value being higher the more the player's payoff varies with his own or with the opponent's action. Consistent with this idea, in applying the model to our experiment (Section 3.3), we maintain that the value of reasoning remains constant within the low-payoff treatments ( $[A]$ ,  $[B]$  and  $[C]$ ), since the payoff structure is identical across these three treatments. Similarly, it is constant within the high-payoff treatments ( $[A+]$ ,  $[B+]$  and  $[C+]$ ). Furthermore, the value of reasoning is higher in these latter treatments, as they have higher payoffs associated with being 'correct'.

**Condition 2** For each  $X = A, B, C$ ,  $v_i^{[X+]}(k) \geq v_i^{[X]}(k)$  for all  $k$ .

It is useful to introduce the following mapping, which identifies the intersection between the value of reasoning and the cost function: Let  $\mathcal{K} : \mathbb{R}_+^{\mathbb{N}} \times \mathbb{R}_+^{\mathbb{N}} \rightarrow \mathbb{N}$  be such that, for any  $(c, v) \in \mathbb{R}_+^{\mathbb{N}} \times \mathbb{R}_+^{\mathbb{N}}$ ,

$$\mathcal{K}(c, v) = \min \{k \in \mathbb{N} : c(k) \leq v(k) \text{ and } c(k+1) > v(k+1)\}, \quad (1)$$

with the understanding that  $\mathcal{K}(c, v) = \infty$  if the set in equation (1) is empty. Player  $i$ 's *cognitive bound*, which represents his understanding of the game, is then determined by the value that this function takes at  $(c_i, v_i)$ :

**Definition 1** Given cost and value functions  $(c_i, v_i)$ , the cognitive bound of player  $i$  is defined as:

$$\hat{k}_i = \mathcal{K}(c_i, v_i). \quad (2)$$

Player  $i$  therefore stops the iterative process when the value of performing an additional round of introspection exceeds the cost. The point at which this occurs identifies his cognitive bound  $\hat{k}_i$ . Note that a player does not compare the benefits and costs at higher  $k$ 's. That is, he does not consider stages of reasoning higher than his current point. A player who has performed  $k$  rounds of introspection is only aware of the portion uncovered by the  $k$  steps, and performs a 'one-step ahead' comparison of the incremental cost and value of reasoning.

This process may appear to translate to a standard optimization problem, in which an agent's marginal cost-marginal benefit analysis can be interpreted as first order conditions of a 'total' value and cost tradeoff. This need not always be the case, however. If the functions

$c_i$  and  $v_i$  cross at most once, then the two procedures are indeed the same. In other words, if the player were able to look ahead to higher  $k$ 's, it would not change his decision to stop the process, since the cost  $c_i$  would remain higher than the benefit  $v_i$  for those  $k$ 's. But suppose that the cost and benefit functions were to cross more than once. Then, with a standard optimization problem, the first crossing would not necessarily be the optimal one. The myopic procedure followed by the players in this model would then not lead to the global optimizer's solution. Thus, player  $i$  does not 'optimize' by taking into account the entire curves; rather, he optimizes locally. The 'myopic' (one-step-ahead) behavior that we assume captures the idea, inherent to the very notion of bounded rationality, that the agents do not know (or are not aware of) what they have not yet thought about. Formalizing this notion is often perceived to be a fundamental difficulty in developing a theory of bounded rationality; in this model it emerges naturally.<sup>16</sup>

Notice that, in our framework, the cognitive bound  $\hat{k}_i$  is monotonic in the level of the cost and of the value functions, irrespective of the shape of these functions:  $\hat{k}_i$  (weakly) decreases as the cognitive costs increase, and it (weakly) increases as the value of reasoning increases.

**Proposition 1** *Under the maintained assumptions of Condition 1:*

1. For any  $c_i$ , the cognitive bound  $\hat{k}_i$  (weakly) increases as  $v_i$  is replaced by  $v'_i$  such that  $v'_i(k) \geq v_i(k)$  for all  $k$ .
2. For any  $v_i$ , the cognitive bound  $\hat{k}_i$  (weakly) decreases as  $c_i$  is replaced by  $c'_i$  such that  $c'_i(k) \geq c_i(k)$  for all  $k$ .

This result is immediate, as shown in the following example of players' reasoning process.

**Example 1** *Figure 1 represents a non-monotonic cost function,  $c_i$ . The value of reasoning,  $v_i$ , is drawn as a constant for simplicity. As in Definition 1, player  $i$ 's cognitive bound,  $\hat{k}_i$ , is determined by the first intersection of the two functions,  $c_i$  and  $v_i$ . In the graph on the left,  $\hat{k}_i = 2$ , meaning that player  $i$  has 'become aware' of one round of reasoning of the opponent. The grey area represents the 'unawareness region' of player  $i$  about the opponents' steps of reasoning of level higher than 1. As the value  $v_i$  is continuously increased (for instance, because the payoffs of the game are increased),  $\hat{k}_i$  remains constant at first, but then increases to  $\hat{k}'_i = 3$  when level  $v'_i$  is reached. Correspondingly, the grey area shifts to the right, uncovering one more round of reasoning of the opponent. If  $v_i$  is further increased,  $i$ 's cognitive bound  $\hat{k}_i$  eventually increases to 4 once  $v^*$  is reached, after which  $\hat{k}_i$  jumps to  $\infty$ . The non-monotonic cost function  $c_i$  thus captures the situation of a player who suddenly understands the game completely after having performed a few rounds of reasoning.*

---

<sup>16</sup>The results we present here, however, would not change if  $\mathcal{K}(\cdot)$  were defined as the optimal solution resulting from a cost-benefit analysis. We maintain the 'myopic' formulation because we find it more natural in the context of level- $k$  reasoning.



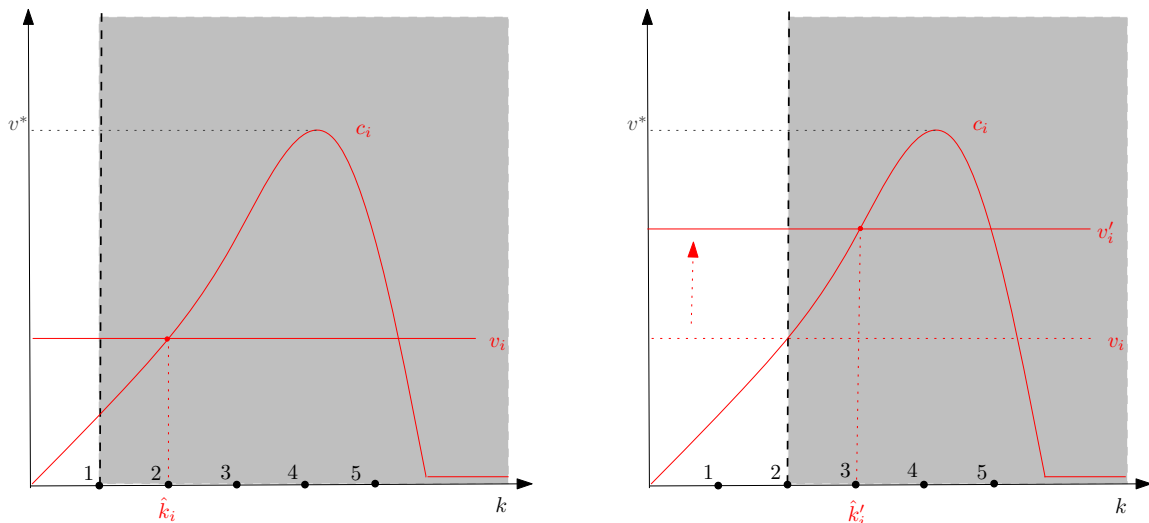


Figure 1: As  $v_i$  increases,  $\hat{k}_i$  (weakly) increases. The grey area represents the ‘unawareness region’ of player  $i$ .

Note that the cost-benefit analysis conducted by the agent and the ensuing bound  $\hat{k}_i$  are independent of the player’s opponent. In this respect, the bound  $\hat{k}_i$  can be seen to be determined separately from the player’s beliefs about the opponent’s sophistication, although both factors affect his behavior.<sup>17</sup>

### 3.2.2 Reasoning about the opponents.

Before choosing an action, players take into account the sophistication of their opponents. Players’ beliefs about their opponents’ sophistication need not be correct, as we do not seek for an equilibrium concept and correctness of beliefs is not guaranteed by introspection alone.

A key conceptual difficulty in level- $k$  models is the notion that a player may face an opponent whom he perceives to be more sophisticated than him. Our model overcomes this limitation and gives meaning to the idea of playing a more sophisticated opponent.

Since a player’s reasoning ability in this model is captured by the cost function, we use the same tool to model players’ beliefs over others’ sophistication. We specifically define sophistication in the following manner:

**Definition 2** Consider two cost functions,  $c'$  and  $c''$ . We say that cost function  $c'$  corresponds to a ‘more sophisticated’ player than  $c''$ , if  $c'(k) \leq c''(k)$  for every  $k$ .

We do not define sophistication directly in terms of the cognitive bounds  $\hat{k}_i$  and  $\hat{k}_j$  because these bounds are a consequence of the cost-benefit analysis, which is determined endogenously by the game. When payoffs are different for the two players, it may be that  $\hat{k}_i < \hat{k}_j$  even if  $i$

<sup>17</sup>Alternatively,  $\hat{k}_i$  can be interpreted as the most sophisticated behavior that player  $i$  would conceive of in the game, if he thinks that his opponent is at least as sophisticated as he is himself. We do not discuss this interpretation of the cognitive bound  $\hat{k}_i$  for the sake of brevity.

is more sophisticated than  $j$ , in the sense of Definition 2. For instance, this may hold if player  $i$  has lower incentives than player  $j$  despite having higher cognitive abilities.

We separate the space of cost functions as follows. For any  $c_i \in \mathbb{R}_+^N$ , let

$$C^+(c_i) = \left\{ c' \in \mathbb{R}_+^N : c_i(k) \geq c'(k) \text{ for every } k \right\} \text{ and}$$

$$C^-(c_i) = \left\{ c' \in \mathbb{R}_+^N : c_i(k) \leq c'(k) \text{ for every } k \right\}.$$

Thus, based on Definition 2,  $C^+(c_i)$  and  $C^-(c_i)$  are comprised of the cost functions that are respectively ‘more’ and ‘less’ sophisticated than  $c_i$ .

**Definition 3** *Let  $c_j^i$  denote player  $i$ 's beliefs over  $j$ 's cost function. Then, player  $i$  ‘believes that his opponent is more (resp., less) sophisticated than himself’, if and only if  $c_j^i \in C^+(c_i)$  (resp.,  $c_j^i \in C^-(c_i)$ ).*

When choosing his action, player  $i$  considers the cost-benefit analysis of his opponent. Given player  $i$ 's beliefs over  $j$ 's costs  $c_j^i$  and given  $j$ 's value function  $v_j$ ,  $i$ 's beliefs about  $j$ 's cost functions would be at the intersection of these two functions, if he is aware of this occurs. Using Definition (1), this point of intersection is  $\mathcal{K}(c_j^i, v_j)$ . However, the maximum bound that he can conceive of for his opponent is constrained by his own cognitive bound. He only conceives of his opponent's cost-benefit analysis in the region uncovered by his own reasoning, which is the region up to  $\hat{k}_i - 1$ . We therefore define player  $i$ 's belief about player  $j$ 's bound,  $\hat{k}_j^i$ , as

$$\hat{k}_j^i = \min \left\{ \hat{k}_i - 1, \mathcal{K}(c_j^i, v_j) \right\}. \quad (3)$$

Equation 3 therefore constrains  $i$ 's beliefs over  $j$ 's bound,  $\hat{k}_j^i$ , to be within the limit of  $i$ 's own understanding. Thus,  $i$ 's beliefs  $c_j^i$  actually represent the *path* of  $i$ 's beliefs about the opponent's reasoning, as  $i$  uncovers more and more steps of reasoning.

We are now in position to specify the behavior of player  $i$ . We assume that player  $i$  believes that  $j$  plays according to his bound as perceived by  $i$ . Letting  $k_j^i$  be  $i$ 's beliefs over the level according to which  $j$  plays, this implies that  $k_j^i = \hat{k}_j^i$ . Player  $i$  best responds to  $k_j^i$ , and therefore plays the action  $a_i^{k_i}$  associated with the ‘behavioral level’,  $k_i$ , defined as

$$k_i = k_j^i + 1. \quad (4)$$

Hence,  $i$  plays according to his own cognitive bound  $\hat{k}_i$  whenever his understanding of  $j$ 's level of play is constrained by his own understanding of the game. If instead player  $i$  believes that he has performed more rounds of introspection than  $j$ , then player  $i$  assumes that  $j$  would play according to *his* ( $j$ 's) maximal bound, and best responds by playing action  $a_i^{k_i}$ .

In brief, there are four  $k$ 's that are determined endogenously. Player  $i$ 's *cognitive bound*,  $\hat{k}_i$ , arises from his cost-benefit analysis, and represents his understanding of the game. Player

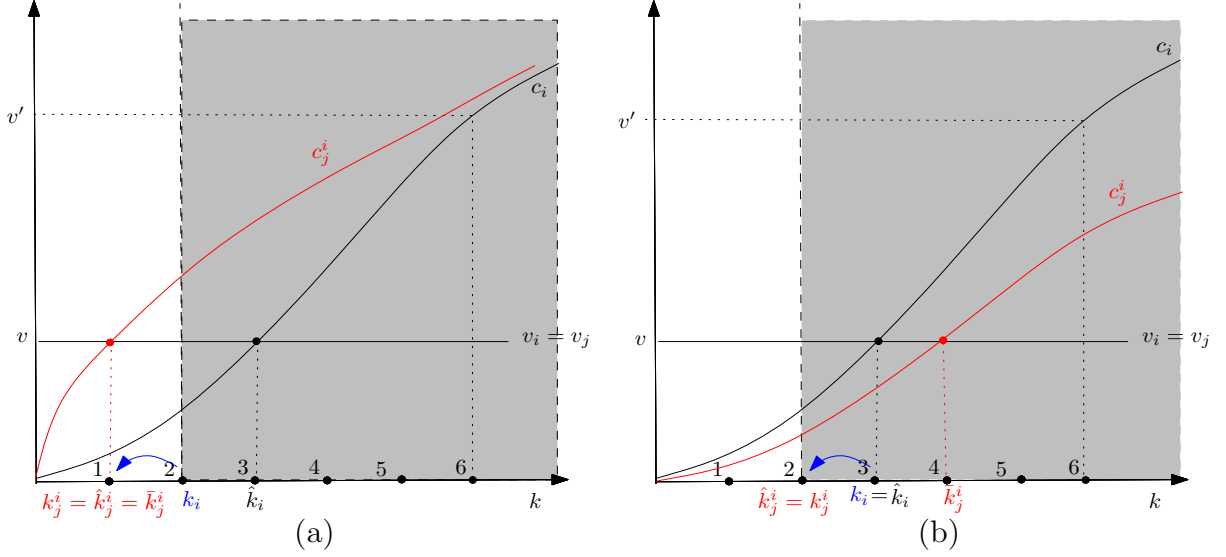


Figure 2: Reasoning about the opponents: on the left,  $c_j^i \in C^-(c_i)$ ; on the right,  $c_j^i \in C^+(c_i)$ . The grey area represents the ‘unawareness region’ of player  $i$ . The intersection of  $c_j^i$  and  $v_j$  is denoted  $\bar{k}_j^i$ .

$i$ 's beliefs about player  $j$ 's understanding of the game,  $\hat{k}_j^i$ , derive from  $i$ 's beliefs about the opponent's cost function ( $c_j^i$ ) and incentives ( $v_j$ ), given  $i$ 's understanding of the game ( $\hat{k}_i$ ). Beliefs  $\hat{k}_j^i$  in turn determine  $i$ 's beliefs about  $j$ 's behavior,  $k_j^i$ . Finally, the ‘behavioral level’ of player  $i$ ,  $k_i$ , follows from  $k_j^i$  and determines the action chosen in the game,  $a_i^{k_i}$ .

**Example 2** Figure 2 represents an agent with cost function  $c_i$  which, given  $v_i$ , induces a cognitive bound of  $\hat{k}_i = 3$ , implying that he has uncovered two rounds of reasoning of his opponent. In Figure 2.a, player  $i$  perceives his opponent to be less sophisticated. Since the intersection between  $c_j^i$  and  $v_j$  falls in the region already uncovered by  $i$ 's cognitive bound,  $i$ 's belief about  $i$ 's cognitive bound is at that point, i.e.  $\hat{k}_j^i = 1$ . This also represents  $i$ 's belief about  $j$ 's behavior,  $k_j^i$ , hence player  $i$  best responds by playing the action associated with level  $k_i = k_j^i + 1 = 2$ . Notice that, in this case, the cognitive bound  $\hat{k}_i$  is not binding:  $k_i < \hat{k}_i$ .

Figure 2.b instead represents the same player reasoning about an opponent that he perceives to be more sophisticated. In this case, the intersection between  $c_j^i$  and  $v_j$  (denoted  $\bar{k}_j^i$  in the graph) falls in the ‘unawareness region’ of player  $i$ . Hence, his perceived cognitive bound for player  $j$  is not  $\bar{k}_j^i$  but  $\hat{k}_j^i = \hat{k}_i - 1$ . Player  $i$  best responds by playing according to level  $k_i = \hat{k}_j^i + 1$ , that is, according to his own cognitive bound  $\hat{k}_i$ . Numerically,  $i$  plays according to  $k_i = 3$  and not  $\bar{k}_j^i + 1 = 5$ .

The proposition below summarizes some implications of our model, which will translate to predictions in the experiment conducted.

**Proposition 2**

1. If  $v_i = v_j$  and if  $c_j^i \in C^+(c_i)$ , then  $k_i = \hat{k}_i$ .
2. If  $v_i = v_j$  and if  $c_j^i \in C^-(c_i)$ , then  $k_i \leq \hat{k}_i$ .
3. If  $v_i = v_j$  and the functions are increased (preserving the symmetry), both  $k_i$  and  $\hat{k}_i$  (weakly) increase.

In words, in a game with symmetric incremental benefits, the cognitive bound  $\hat{k}_i$  is always binding for a player who believes that he is playing against a more sophisticated opponent. If instead he perceives his opponent to be less sophisticated, then his cognitive bound  $\hat{k}_i$  is lower than his behavioral  $k_i$ . Taken together, these two points imply that, in the modified 11-20 game, players play (weakly) higher actions against opponents that they perceive to be less sophisticated than against those that they perceive to be more sophisticated. Lastly, if all players' value of reasoning increase then  $i$ 's cognitive bound  $\hat{k}_i$  and his behavioral  $k_i$  both increase. In the modified 11-20 game, this point implies that players play lower actions as payoffs are increased.

**3.2.3 Higher Order Reasoning**

The analysis above is based on the implicit assumption that players believe that their opponents have correct beliefs about their own cost function. This assumption is used to conclude that, if player  $i$  believes that his opponent has performed less rounds of introspection ( $\hat{k}_j^i < \hat{k}_i$ ), then his opponent would play according to his own cognitive bound ( $k_j^i = \hat{k}_j^i$ ). The model however can be extended to accommodate higher order uncertainty, allowing players to view their opponents' beliefs to be incorrect. This extension is relevant for the analysis of treatments  $[C]$  and  $[C+]$ , relative to  $[B]$  and  $[B+]$ , in the experiment. Specifically, let  $c_i^{ij}$  denote  $i$ 's beliefs about  $j$ 's beliefs about  $i$ 's cost function. The model discussed so far implicitly assumes that  $c_i^{ij} = c_i$ . This, however, need not be the case if player  $i$ 's opponent plays against someone other than player  $i$ . In treatment  $[C]$  for instance, a player faces the action of an opponent playing someone else. In that case, we explicitly model  $i$ 's beliefs about  $j$ 's opponent.

Let the triple  $(c_i, c_j^i, c_i^{ij})$  represent player  $i$ 's cost  $c_i$ , his beliefs over his opponents cost  $c_j^i$ , and his beliefs over his opponent's beliefs about his own cost,  $c_i^{ij}$ . We extend the model by having player  $i$  take into account that player  $j$  may play at a lower  $k$  than his cognitive bound. In particular, while  $\hat{k}_j^i$  still defines player  $i$ 's perception of player  $j$ 's cognitive bound, it does not define his perception of player  $j$ 's actual play. Instead, player  $i$  considers player  $j$ 's reasoning about  $i$  using the pair  $(c_j^i, c_i^{ij})$ . That is, he puts himself in player  $j$ 's 'shoes' and calculates  $j$ 's opponent's bound from his standpoint.<sup>18</sup> Formally, player  $i$  views player  $j$ 's perception of player  $i$ 's limit to be

$$\hat{k}_i^{ij} = \min \left\{ \mathcal{K} \left( c_i^{ij}, v_i \right), \hat{k}_j^i - 1 \right\}. \tag{5}$$

---

<sup>18</sup>More precisely, he uses the triple  $(c_j^i, c_i^{ij}, c_j^{ji})$  for his opponent, where  $c_j^{ji} = c_j^i$ . Hence, the implicit assumption is that  $i$ 's beliefs of order higher than 2 are consistent with common certainty of  $c_j^i$ .

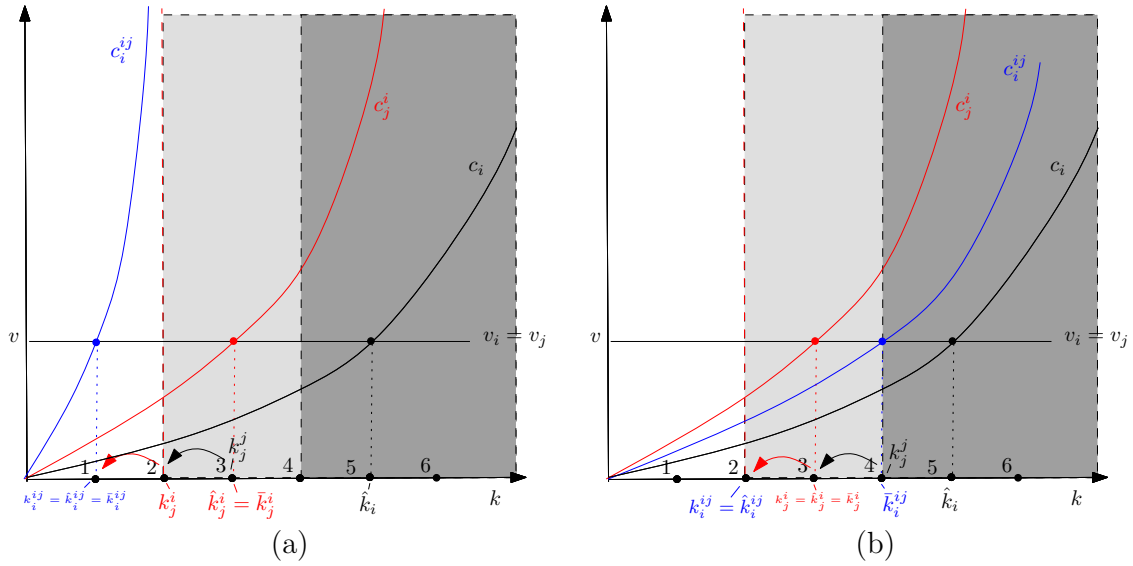


Figure 3: Higher Order Reasoning:  $c_j^i \in C^-(c_i)$ , with  $c_i^{ij} \in C^-(c_j^i)$  on the left, and with  $c_i^{ij} \in C^+(c_j^i)$  on the right. The dark grey area represents the ‘unawareness region’ of player  $i$ , whose cognitive bound is  $\hat{k}_i = 5$ . The light grey area represents the unawareness region of  $j$ , as perceived by  $i$ . The intersection of  $c_j^i$  and  $v_j$  is denoted  $\hat{k}_j^i$ , and the intersection of  $c_i^{ij}$  and  $v_j$  is denoted  $\bar{k}_i^{ij}$ .

If player  $i$  is capable of conceiving of such a level, that is, if  $\hat{k}_i^{ij} + 1 \leq \hat{k}_i - 1$ , then he expects  $j$  to play according to level  $\hat{k}_i^{ij} + 1$ . Hence, we define  $i$ ’s perception of  $j$ ’s play as

$$k_j^i = \min \left\{ \hat{k}_i^{ij} + 1, \hat{k}_i - 1 \right\}. \quad (6)$$

Player  $i$  then best responds playing action  $a_i^{k_i}$ , where  $k_i = k_j^i + 1$ .<sup>19</sup>

**Example 3** Figure 3 represents a player with cost function  $c_i$  reasoning about an opponent that he regards as less sophisticated. The dark grey area represents the ‘unawareness region’ of player  $i$ , whose cognitive bound is  $\hat{k}_i = 5$ . The light grey area represents the ‘unawareness region’ of  $j$ , as perceived by  $i$ .

In Figure 3.a, player  $i$  believes that  $j$  thinks that  $i$  is even less sophisticated (that is,  $c_i^{ij} \in C^-(c_j^i)$ ). Given beliefs  $c_j^i$ ,  $i$ ’s perception of  $j$ ’s cognitive bound is  $\hat{k}_j^i = 3$ . This bound, however, does not necessarily coincide with  $k_j^i$ , which is determined by the intersection of  $c_j^i$  and  $v_i$ , which occurs at 1. This falls in the region that  $i$  perceives  $j$  to have uncovered. Hence,  $\hat{k}_i^{ij} = k_i^{ij} = 1$  and  $i$  believes that  $j$ ’s best response is to play according to  $k_j^i = 2$ . In turn,  $i$ ’s best responds by playing according to  $k_i = 3$ .

In Figure 3.b,  $c_i^{ij} \in C^+(c_j^i)$ , and therefore  $i$  believes that  $j$  views him as more sophisticated.

<sup>19</sup>Note that setting  $c_i^{ij} = c_i$ , we obtain the case of the previous subsection, where  $k_j^i = \hat{k}_j^i$ . This is so because, in that case,  $\mathcal{K}(c_i^{ij}, v_i) = \hat{k}_i$ , hence (5) delivers  $\hat{k}_i^{ij} = \hat{k}_j^i - 1$ . By definition of  $\hat{k}_j^i$  (3),  $\hat{k}_j^i - 1 < \hat{k}_i - 1$ , hence by (6),  $k_j^i = \hat{k}_i^{ij} + 1 = \hat{k}_j^i$ .

Player  $i$  therefore expects  $j$  to play at his maximum bound,  $k_j^i = \hat{k}_j^i = 3$ .<sup>20</sup> The best response is thus to play according to  $k_i = 4$ .

### 3.3 Theoretical Predictions for the Experiment

Recall that we use the terminology ‘label  $I$ ’ (resp., ‘label  $II$ ’) to refer indiscriminately to the ‘high score’ (‘low score’) subjects in the endogenous classification, or to the ‘math and sciences’ (‘humanities’) subjects for the exogenous. Accordingly, we introduce notation  $l_i = \{I, II\}$  to refer to player  $i$ ’s label. We describe next how we apply our theory to the treatments of the experiment.

The assumption that the value of reasoning is determined solely by players’ payoffs implies that  $v_i$  remains constant throughout the low-payoff treatments  $[A], [B], [C]$ , as well as throughout the high-payoff treatments  $[A+], [B+], [C+]$ . In addition, Condition 2 (p. 14) requires that  $v_i$  increases when moving from the low to the high-payoff treatments. We also assume that an agent’s cognitive ability,  $c_i$ , remains constant throughout all treatments. His beliefs,  $c_j^i$  and  $c_i^{ij}$ , however, change with his opponent’s label, although not with the game payoffs. Hence, for  $X = A, B, C$ ,  $c_j^{i,[X]} = c_j^{i,[X+]}$  and  $c_i^{ij,[X]} = c_i^{ij,[X+]}$ .

To apply our model to the experiment we posit that whenever a subject plays against a label  $I$  opponent (resp., label  $II$ ), then he believes that the opponent is more (less) sophisticated than himself. For higher order beliefs, we assume that  $c_i^{ij,[X]} = c_i$  for  $X = A, B$ . In treatment  $[C]$ , we apply to  $j$ ’s beliefs the same logic applied to  $i$ ’s first order beliefs: if player  $j$ ’s opponent has label  $I$  (resp.,  $II$ ), then  $c_j^{ij} \in C^+(c_j^i)$  (resp.,  $c_j^{ij} \in C^-(c_j^i)$ ). Hence,  $c_i^{ij,[C]} \in C^-(c_j^i)$  if  $l_i = I$  and  $c_i^{ij,[C]} \in C^+(c_j^i)$  if  $l_i = II$ . Define first a player’s *type* to be:

$$t_i = \left( c_i, c_j^{i,I}, c_j^{i,II}, c_i^{ij,[C]} \right),$$

where  $c_i$  is his cost function;  $c_j^{i,I}$  are his beliefs when he plays against a label  $I$  subject;  $c_j^{i,II}$  his beliefs when he plays against a label  $II$  subject;  $c_i^{ij,[C]}$  are his higher order beliefs in treatments  $[C]$  and  $[C+]$ .

Define the sets  $T^*$ ,  $T_I^*$  and  $T_{II}^*$  as follows:

$$T^* = \left\{ t_i \in \left( \mathbb{R}_+^{\mathbb{N}} \right)^4 : c_j^{i,I} \in C^+(c_i) \text{ and } c_j^{i,II} \in C^-(c_i) \right\},$$

$$T_I^* = \left\{ t_i \in T^* : c_i^{ij,[C]} \in C^-(c_j^{i,II}) \right\}, \quad (7)$$

$$\text{and } T_{II}^* = \left\{ t_i \in T^* : c_i^{ij,[C]} \in C^+(c_j^{i,I}) \right\}. \quad (8)$$

In words,  $T^*$  denotes the set of all possible types  $t_i = \left( c_i, c_j^{i,I}, c_j^{i,II}, c_i^{ij,[C]} \right)$ , whereas the sets  $T_I^*$  and  $T_{II}^*$  denote the set of types for label  $I$  and  $II$ , respectively, that are consistent with the assumptions discussed above. This is summarized by the following conditions:

<sup>20</sup>This example shows that the result from the previous subsection that  $k_j^i = \hat{k}_j^i$  if  $c_j^i \in C^-(c_i)$ , requires only that  $c_i^{ij} \in C^+(c_j^i)$ ; it is not necessary that  $c_i^{ij} = c_i$ .

**Condition 3** If  $i$  is a subject with label  $l_i$ , then  $t_i \in T_{l_i}^*$ , where  $l_i \in \{I, II\}$ .

The next proposition summarizes the implications of the model under these assumptions.

**Proposition 3** For any treatment  $[X]$ , let  $a_i^{[X]} \in \{11, 12, \dots, 20\}$  denote  $i$ 's action in treatment  $[X]$ . Under Conditions 1, 2 and 3, the following holds: (1) For any player  $i$  and for any  $X = A, B, C$ ,  $a_i^{[X+]} \leq a_i^{[X]}$ ; (2) If  $l_i = I$ , then  $a_i^{[A]} \leq a_i^{[B]} \leq a_i^{[C]}$  and  $a_i^{[A+]} \leq a_i^{[B+]} \leq a_i^{[C+]}$ ; if  $l_i = II$ , then  $a_i^{[A]} \geq a_i^{[B]} = a_i^{[C]}$  and  $a_i^{[A+]} \geq a_i^{[B+]} = a_i^{[C+]}$ .

Applying Proposition 3 leads to unambiguous first order stochastic dominance relations between the distributions of actions in the different treatments.<sup>21</sup> Specifically, suppose that the subjects in our experiment of label  $I$  are independently drawn from a distribution  $G_I$  over  $T_I^*$ , and subjects from label  $II$  are independently drawn from a distribution  $G_{II}$  over  $T_{II}^*$ . Let  $F_X^l$  denote the cumulative distribution of actions  $a \in \{11, \dots, 20\}$  in treatment  $X$  for label  $l \in \{I, II\}$ . Then, letting  $\succsim$  denote the first order stochastic dominance relation, it is immediate from Proposition 3 that  $F_C^I \succsim F_B^I \succsim F_A^I$ , and that  $F_A^{II} \succsim F_B^{II}$  and  $F_A^{II} \succsim F_C^{II}$ . Distributions  $F_B^{II}$  and  $F_C^{II}$ , however, should coincide. Intuitively, a label  $I$  subject plays according to a (weakly) lower action (higher  $k$ ) in treatment [B] than in [A], since he believes that he is playing a less sophisticated player. He plays an even lower action against a less sophisticated opponent who himself believes he is playing a less sophisticated player (as in [C]). Similarly, a label  $II$  subject plays according to a higher action in treatment [B] than [A]. Given that this subject plays according to his cognitive bound in [B], he will still play according to this bound in [C]. The same results hold when comparing treatments [A+], [B+] and [C+]. Lastly, for all  $X = A, B, C$  and  $l = I, II$ , we have that  $F_X^l \succsim F_{X+}^l$ : as payoffs are increased, players perform more rounds of reasoning (and they expect their opponents to do so), hence they play lower actions.

These results, summarized in Table 2, remain true under several variations of the assumptions entailed by the definition of the sets  $T_I^*$  and  $T_{II}^*$ . For instance, the results hold if we allow subjects to be ‘mis-labeled’, in the sense that label  $I$  ( $II$ ) subjects may believe that they are less (more) sophisticated than label  $II$  ( $I$ ) subjects, provided that they kept the belief that label  $I$  subjects are known to be more sophisticated than label  $II$  subjects. The framework can also be extended to allow for a stochastic component to the reasoning process, which would introduce an element of noise to the stochastic dominance properties described above.

### 3.4 Related Models

The main feature of our framework consists of modeling the cognitive bound as a consequence of a tradeoff between value of reasoning and costs of cognition. The general notion that players follow a cost-benefit analysis is present in the language of Camerer, Ho and Chong (2004), but not in the cognitive hierarchy (CH) model itself, as players’ cognitive types remain exogenous. A recent paper by Choi (2012) extends the CH model by letting cognitive types result from an

<sup>21</sup>Given two cumulative distributions  $F(x)$  and  $G(x)$ , we say that  $F$  (weakly) first order stochastically dominates  $G$  if  $F(x) \leq G(x)$  for every  $x$ .

Labels	Changing beliefs (low payoffs)	Changing beliefs (high payoffs)	Changing payoffs
$l_i = I$	$F_C \succsim F_B \succsim F_A$	$F_{C+} \succsim F_{B+} \succsim F_{A+}$	$F_X \succsim F_{X+}$ for $X = A, B, C$
$l_i = II$	$F_B \succsim F_A; F_B \approx F_C$	$F_{A+} \succsim F_{B+}; F_{B+} \approx F_{C+}$	$F_X \succsim F_{X+}$ for $X = A, B, C$

Table 2: Summary of the theoretical predictions of the relations between the distribution of actions in different treatments.

optimal choice. This optimization serves to provide identification restrictions to estimate the distribution of types across different environments, and is motivated through an evolutionary argument. In contrast, our goal is to provide an explicit model of reasoning, in which players think about both the game and about the reasoning process of their opponents. The objectives and modeling choices are therefore distant.

Gabaix (2012) also proposes a framework in which the accuracy of players’ beliefs about the opponents’ behavior is determined by a cost-benefit analysis. At a conceptual level, the main goal of Gabaix (2012) is to provide an equilibrium concept that allows players to have both incorrect beliefs as well as to respond non-optimally. Our model focuses on agents’ cognitive limitations in reasoning about the opponents, and ignores the orthogonal issue of limitations in computing best responses. Furthermore, unlike Gabaix (2012) and CH models of Camerer, Ho and Chong (2004) and Choi (2012), our framework does not impose equilibrium-like restrictions that relate agents’ beliefs to the actual distribution of actions or types. Since our model is purely one of introspection, we avoid conditions that would not be a consequence of players’ reasoning process alone. This feature of our model is shared by the level- $k$  models of Nagel (1995) and Crawford and Iriberry (2007).<sup>22</sup>

## 4 Experimental results

We present, for brevity, only the experimental results for the grouped exogenous and endogenous classifications.<sup>23</sup> We pool the label  $I$  subjects (‘math and sciences’ for exogenous treatments and ‘high’ for endogenous treatments), and we pool the label  $II$  subjects (‘humanities’ for exogenous treatments and ‘low’ for endogenous treatments). Recall that we do not take these labels to indicate actual game theoretical sophistication, but of *perceived* sophistication by the subjects. Moreover, we present the results by pooling together the treatments when they are repeated. For these repetitions, our pooling is justified by tests for equality of distribution. We analyze first the results when subjects’ payoffs are changed, followed by the results when their beliefs over opponents are varied. Overall, the results appear to confirm the theory. We consider different statistical tests, including Mann-Whitney-Wilcoxon (rank sum)

<sup>22</sup>The ‘noisy introspection’ model of Goeree and Holt (1999) introduces noisy responses in a non-equilibrium solution concept.

<sup>23</sup>The figures for the separate classifications are consistent with the results for the grouped classifications. They are available upon request, as are other supplementary materials.



and Kolmogorov-Smirnov tests. To further test our first order stochastic dominance relations, we use the method introduced by Davidson and Duclos (2000) (henceforth, DD). We also conduct regressions, and find that most of the relevant coefficients are significant with the sign expected. This lends support to our model, and to the general claim that beliefs and incentives affect the individuals' depth of reasoning and their behavior.

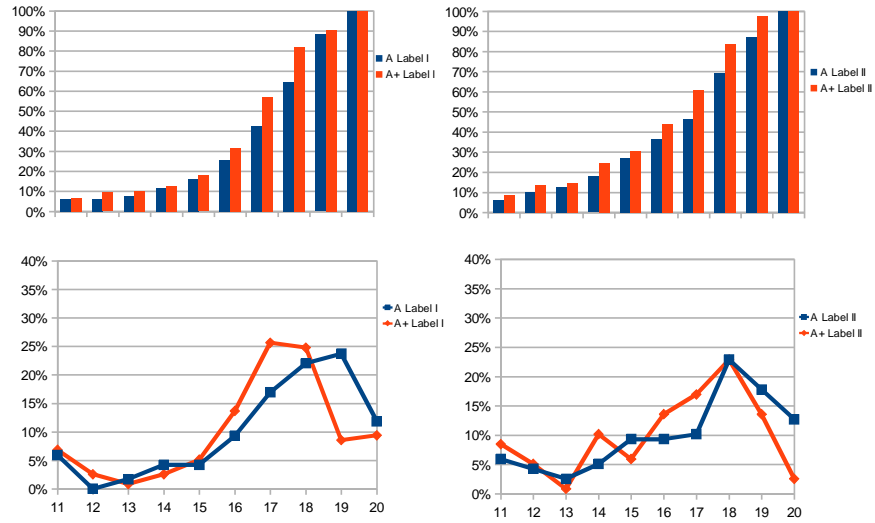
In the (random-effects) ordinary least squares estimations (OLS) that follow, we regress, for each label, the outcome on a dummy for the treatments, and another for the classification (endogenous or exogenous). The latter is never significant. To control for 'feedback-free' learning, we exploit two factors. First, we use randomization of treatments, particularly within [A], [B] and [C], and within [A+], [B+] and [C+]. Second, we exploit the repetition of treatments to do equality of distribution tests for the same treatment. For instance, each agent plays treatment [A] twice with other treatments in between, and we test for equality of distribution between the first and the second time treatment [A] occurs. The results of Mann-Whitney-Wilcoxon tests and Kolmogorov-Smirnov tests are highly suggestive that the distributions are equal. As an added robustness check in comparing low to high payoffs in the OLS estimations, we also control for the round at which they are played, and find that it is never significant.

All regressions and statistical tests are in Appendix C. The OLS regressions are in Table 5 of Appendix C and the Mann-Whitney-Wilcoxon and Kolmogorov-Smirnov tests for changes in payoffs and beliefs over opponents are in Table 6 and Table 7, respectively. The DD tests are in Tables 8 and 9. The equality of distributions tests for the order of treatments are in Table 10.

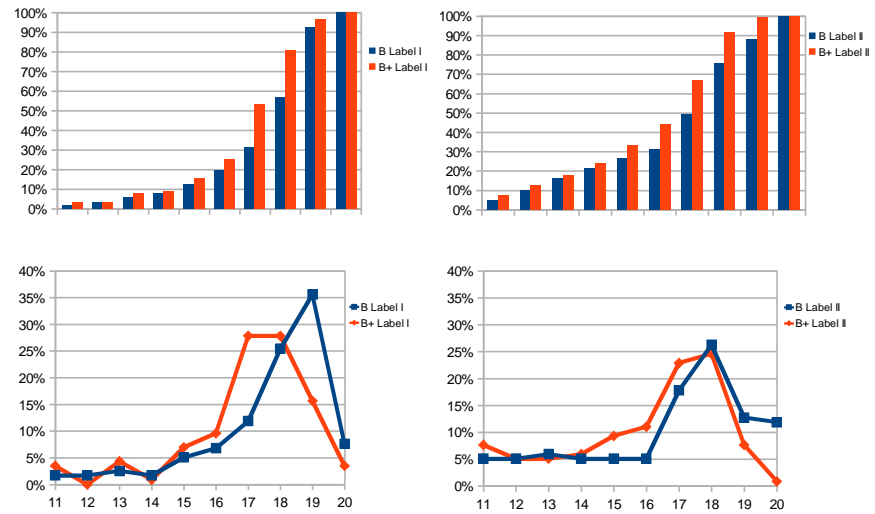
#### 4.1 Changing incentives

As the value of introspection increases for players and their opponents, the model predicts that they would choose actions associated with higher  $k$ 's. Specifically, comparing treatments across different marginal values of payoffs,  $F_A \succeq F_{A+}$ ,  $F_B \succeq F_{B+}$  and  $F_C \succeq F_{C+}$ . These implications hold for both label *I* and label *II* subjects. Beginning with label *I*, it is clear from Figure 4.a (left) that the empirical distribution [A] stochastically dominates [A+] everywhere. Furthermore, distribution [B] stochastically dominates [B+] everywhere, and [C] clearly stochastically dominates [C+] everywhere (Figures 4.b and 4.c). Using a DD test for each comparison, we find that these results are indeed consistent with first order stochastic dominance, as shown by all the signs of the statistics. These results are therefore consistent with our theoretical predictions. Conducting an OLS regression, we find that the coefficients are highly significant ( $< 1\%$ ) for distributions [A] compared to [A+], [B] to [B+] and [C] to [C+], and of the correct sign. The Mann-Whitney-Wilcoxon tests and Kolmogorov-Smirnov tests are both significant ( $< 5\%$ ) or highly significant ( $< 1\%$ ) for all of these comparisons of distribution as well.

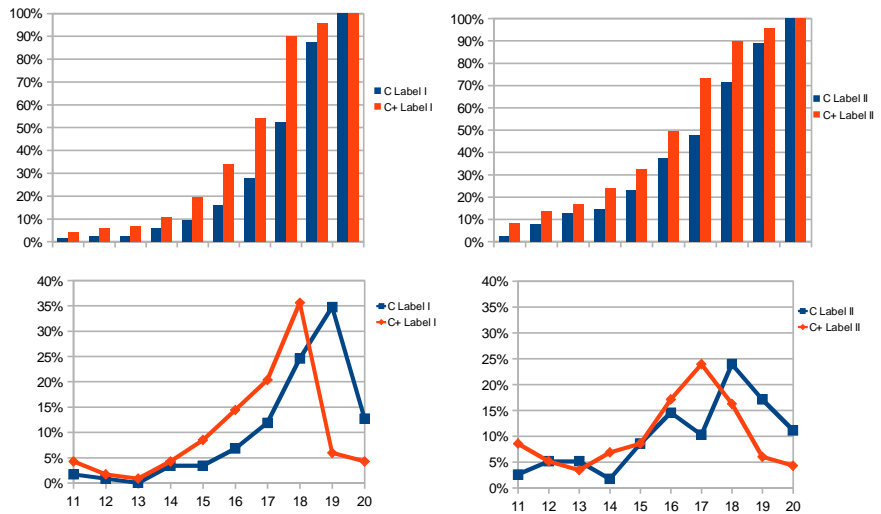
These findings are consistent with the model, and with the view that agents perform more rounds of reasoning if the incentives are increased. These results also indicate that changing from an extra 20 tokens to an extra 80 tokens determines a large enough shift in the value



(a) : A and A+



(b) : B and B+



(c) : C and C+

Figure 4: Changing Payoffs, label *I* (left) and label *II* (right); cumulative distributions (top) and frequency distributions (bottom)

function that it leads agents to increase their level of reasoning.<sup>24</sup> The graphs in Figure 4 are suggestive of the idea that there is a shift in the distribution. We note as well that level-1 and level-2 play is modal for treatments [A], [B] and [C], while level-2 and level-3 play is modal for [A+], [B+] and [C+]. The means of the distributions change from 17.3, 17.7 and 18.0 for treatments [A], [B] and [C], respectively, to 16.8, 17.1 and 16.9 for treatments [A+], [B+] and [C+], respectively.

For label *II*, the stochastic dominance relationships hold everywhere for all three comparisons, [A] to [A+], [B] to [B+] and [C] to [C+], as shown in Figure 4 (right). All the signs of the statistics of the DD tests also confirm that these results are consistent with stochastic dominance, and the coefficients from the OLS regression are significant ( $< 5\%$ ) for distributions [A] compared to [A+] and highly significant ( $< 1\%$ ) for distributions [B] compared to [B+] and [C] to [C+]. These coefficients are of the correct sign. The Mann-Whitney-Wilcoxon test is significant ( $< 5\%$ ) or highly significant ( $< 1\%$ ) for all of these comparisons, and the Kolmogorov-Smirnov test is significant ( $< 5\%$ ) for [B] compared to [B+] and highly significant ( $< 1\%$ ) for [C] compared to [C+]. As with label *I*, the figures for label *II* are suggestive of a shift in the distribution. Here as well level-1 and level-2 play is modal for treatments [A] and [C] (level-2 and level-3 play is modal for [B]), while level-2 and level-3 play is modal for [A+],[B+] and [C+]. The means of the distributions change from 16.9, 16.8 and 16.9 for treatments [A], [B] and [C], respectively, to 15.6, 15.6 and 15.5 for treatments [A+], [B+] and [C+], respectively.

## 4.2 Changing beliefs about the opponents

Consider the comparison between homogeneous treatment [A], heterogeneous treatment [B] and replacement treatment [C]. According to the theoretical model,  $F_C \succsim F_B \succsim F_A$  for label *I* players. This result seems consistent with the data displayed in Figure 5. Distribution [C] clearly stochastically dominates [B] everywhere, and [B] stochastically dominates [A] nearly everywhere.<sup>25</sup> We also note that [C] clearly stochastically dominates [A] everywhere.

Note that the theoretical model does not predict *strict* stochastic dominance relations  $F_C \succ F_B$  or  $F_B \succ F_A$ . That is, it allows for individuals' beliefs over their relative sophistication to be such that they would play the same against lower or higher levels of sophistication. The distinct pattern that emerges from Figure 5 indicates that label *I* individuals view the cost function associated with label *II* as sufficiently far from their own to induce a difference in their chosen action. We also note that this pattern could not be explained by existing models of level- $k$  reasoning, as they do not endogenize that level of play may vary even if the payoffs of the game remain constant.

The OLS estimates comparing [A] to [B] are significant at  $< 10\%$ , and the estimates comparing [A] to [C] are highly significant ( $< 1\%$ ). The estimates comparing [B] to [C],

<sup>24</sup>See the discussion in Section 5.2 on disentangling  $\hat{k}_i$  from  $k_j^i$ .

<sup>25</sup>The only exception is at action 19, which is consistent with the well-known observation that stochastic dominance relations are often violated near the endpoints, even when the true distributions are stochastically-dominance ranked.

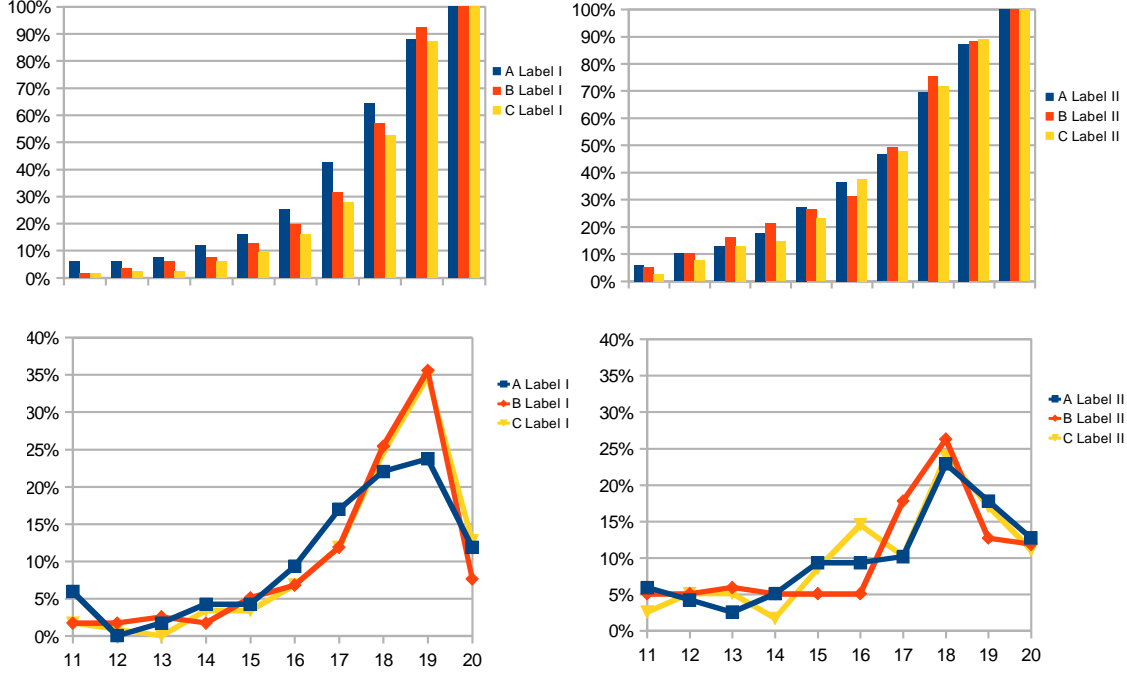


Figure 5: Treatments [A], [B] and [C] for label *I* (left) and *II* (right); cumulative distributions (top) and frequency distributions (bottom)

however, are not significant. Figure 5 reveals that distributions [B] and [C] remain very close to each other, and so the lack of significance is not surprising.

Turning next to label *II* players, the model predicts  $F_A \succsim F_B \approx F_C$ . Here, no clear difference emerges from Figure 5 between the three cumulative distributions. Conducting Mann-Whitney-Wilcoxon and Kolmogorov-Smirnov equality of distribution tests confirms the visual intuition, and the OLS estimates are not significant for any of the comparisons of [A] to [B], [B] to [C] or [A] to [C]. While  $F_B \approx F_C$  is the exact prediction of the theoretical model, the result that  $F_A \approx F_B$  indicates that label *II* subjects do not view the sophistication of other label *II* subjects as significantly lower than their own, and therefore do not adjust their level of play in a measurable way. This result can therefore serve as a first step towards identifying subjects' beliefs over their opponents.

We now compare the high payoff treatments [A+], [B+] and [C+], in which an individual whose action is exactly one below his opponent's receives an additional 80 tokens rather than 20. The model makes analogous predictions for these cases as it does for treatments [A], [B] and [C], namely that  $F_{C+} \succsim F_{B+} \succsim F_{A+}$  for label *I* and  $F_{C+} \succsim F_{B+} \approx F_{A+}$  for label *II*. From Figure 6, no discernible pattern emerges either for label *I* or for label *II*, and we note that the (frequency) distributions are close to each other. Both Mann-Whitney-Wilcoxon and Kolmogorov-Smirnov tests for equality of distribution are consistent with this reading. None of the OLS estimates for the comparison of [A+] and [B+] or [A+] and [C+] are significant, for either label *I* or label *II*. The OLS estimates of [B+] and [C+] for label *II* are not significant

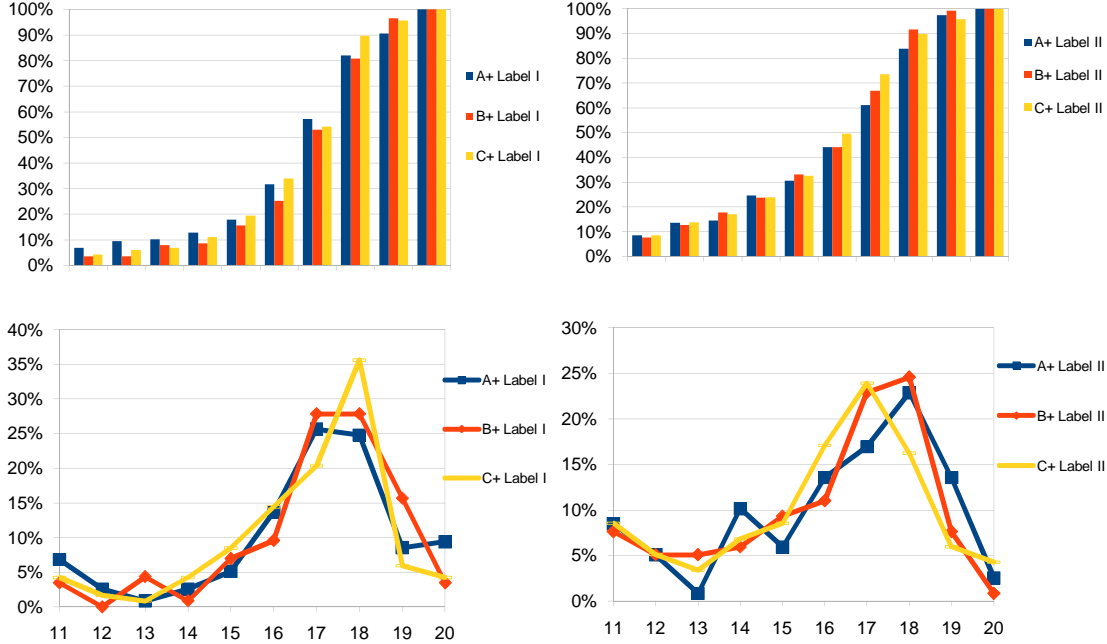


Figure 6: Treatments [A+], [B+] and [C+] for label *I* (left) and *II* (right); cumulative distributions (top) and frequency distributions (bottom)

either, although we note that they are significant from [B+] to [C+] ( $< 5\%$ ) for label *I* in the other direction. Given the closeness of the distribution as confirmed by the equality of distribution tests, this significance does not seem to be of first order.

The last results comparing treatments [A+], [B+] and [C+], viewed together with the results for treatments [A], [B] and [C], are indicative of label *I* subjects' beliefs. Specifically, these results suggest that label *I* subjects believe that the cost functions associated with label *II* subjects are higher than their own at low levels of  $k$ , but become closer to their own cost function at higher  $k$ 's. In other words, label *I* subjects believe that, when sufficiently motivated, label *II* subjects are essentially the same as label *I*. An example of cost functions that satisfy this property is provided in Figure 2.a (p. 18). While the present analysis does not allow for precise identification of subjects' cost function, an extension of our approach could be used for this purpose.

## 5 Additional Treatments

We now describe additional treatments conducted on a subset of the subjects with the aim of exploring less immediate predictions of our theory as well as possible directions for future research. These designs are more complex and cognitively demanding for the subjects, and some rely on the theoretical model more than they directly test it. Overall, the results are encouraging, especially in light of these factors. The instructions of these additional treatments

Treatment	Own label	Opponent's label	Own payoffs	Opponent's payoffs	Own Tutorial	Opponent's Tutorial
Tutorial [D]	$I (II)$	$I (II)$	Low	Low	Yes	Yes
Mixed tutorial [E]	$I (II)$	$II (I)$	Low	Low	Yes	No
Mixed tutorial [F]	$I (II)$	$II (I)$	Low	Low	Yes	No

Table 3: Treatment summary for post-tutorial treatments [D], [E] and [F]

are in Appendix A.1.3 (treatments [D], [E] and [F]) and Appendix A.1.4 (treatments [K] and [L]).

## 5.1 Identifying Beliefs

The theoretical model of Section 3 offers a clear distinction between a player's cognitive bound,  $\hat{k}_i$ , and his behavioral  $k_i$ , which is determined by his beliefs about the opponent,  $k_j^i$ . Recall that we have made the simplifying assumption that subjects view label  $I$  opponents to be more sophisticated than they are themselves. As previously mentioned, this condition is not necessary for our theoretical predictions. Here, we present a test of whether players overall do indeed play according to their bound  $\hat{k}$  when playing label  $I$  opponents. As we demonstrate below, the evidence seems consistent with this assumption.

Specifically, we aim to test whether subjects play according to their own cognitive bound in treatments [A] and [B] or whether they are responding to their beliefs about the opponents' cognitive bound. In the context of our model, this requires a mechanism for setting the players' own costs to zero while holding their beliefs about their opponents constant. For instance, consider the example of Figure 2 (p. 18). If we change the cost function of player  $i$  so that  $c_i(k) = 0$  for every  $k$ , then player  $i$  would play according to  $k_j^i = \bar{k}_j^i$ . If the cognitive bound had not been binding before (as in Figure 2.a), then the agent's behavior would remain the same. If instead his cognitive bound had been binding beforehand (Figure 2.b), then his behavior may change, since  $k_j^i < \bar{k}_j^i$  before the decrease in cost. Hence, within our model, the choice of a lower number (higher  $k$ ) suggests that the player's action had been determined by his cognitive bound rather than his beliefs. Our treatments [D], [E] and [F], summarized in Table 5.1, are designed precisely to operationalize this thought experiment.

After having administered the main treatments, we expose all eighty subjects from four of the six sessions (two for the endogenous and two for the exogenous classifications) to a 'game theory tutorial'. This tutorial explains how, through the chain of best replies, 'infinitely sophisticated and rational players' would play (11, 11). We interpret offering this explanation as setting the subjects' cognitive costs to zero. We then proceed with three new (post-tutorial) treatments, each repeated twice. In treatment [D], we instruct each subject to play the baseline game (with low payoffs) against another subject who has also been given the same tutorial. Not surprisingly, a high fraction of the subjects (48% of label  $I$  and 55% of label  $II$ ) announce 11. In treatment [E], we instruct the subjects who had previously received the tutorial to play the baseline game against a player of the same label who had *not* received the tutorial (that is, as in the homogenous treatment [A]). Analogously, treatment [F] contains the same

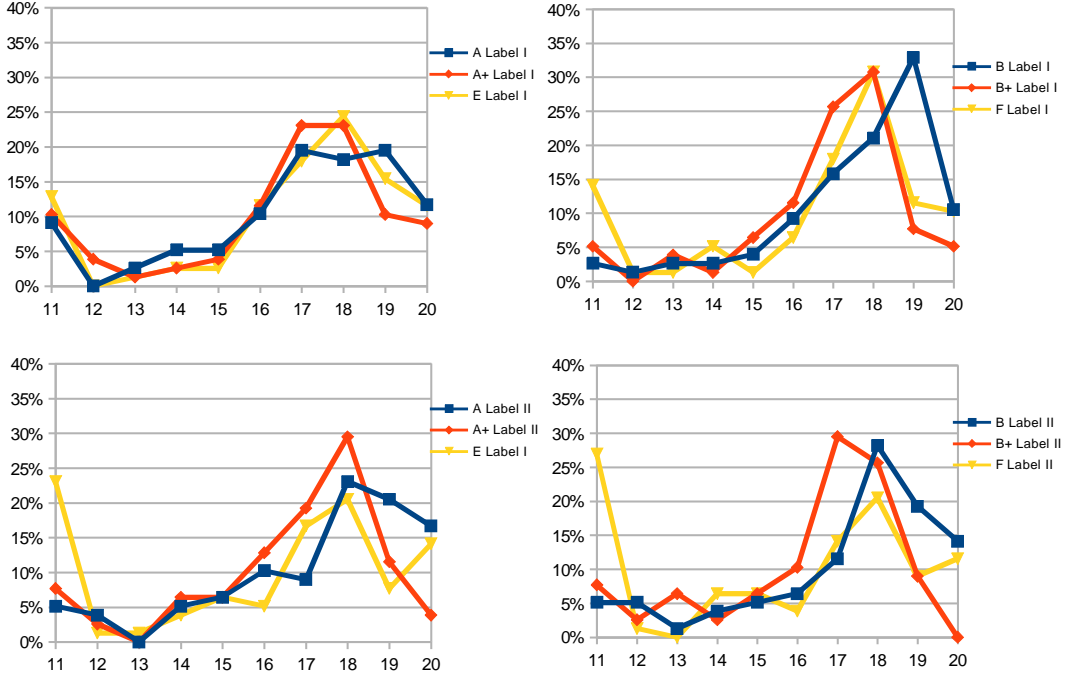


Figure 7: Post-tutorial treatment [E] compared to pre-tutorial treatments [A] and [A+]; post-tutorial treatment [F] compared to [B] and [B+] for labels *I* and *II*; frequency distributions

instructions but with the subjects facing an opponent from a different label (as in [B]). Hence, subjects essentially face the same opponents in treatments [E] and [A] (and in [F] and [B]), but their cost function has been ‘turned off’ in treatments [E] and [F]. If their cognitive bound are *not* binding in (pre-tutorial) treatments [A] and [B], then the distributions of actions in (post-tutorial) treatments [E] and [F] should be the same as in [A] and [B], respectively. Figure 7 displays the results for the two labels. Comparing [A] to [E] and [B] to [F], we observe that the distributions of actions shift to the left, with different degrees across labels and treatments. Through the lens of our model, the data suggest that the cognitive bound for label *I* in treatment [A] and for label *II* in treatment [B] are binding for at least some of the agents.

We emphasize, however, that this interpretation should be taken with a grain of salt. It is not obvious that a highly manipulative intervention such as providing the tutorial would change players’ cost function without affecting their beliefs. Nonetheless, given the complexity of the design and the instructions, we find these results to be encouraging.<sup>26</sup>

Figure 7 also shows the distributions for the treatments with high payoffs, [A+] and [B+]. This figure reveals how, as predicted by our model, increasing the incentives (from [A] to [A+] and from [B] to [B+]) produces effects analogous to reducing the cost function (from [A] to [E] and from [B] to [F]). Unlike treatments [E] and [F], however, subjects in treatments [A+]

<sup>26</sup>Interestingly, a relatively high percentage of subjects from label *II* play 11 in post-tutorial treatments [E] and [F], compared to [A] and [A+]. Furthermore, the percentage of label *II*’s who play 11 in [E] and [F] is nearly twice the percentage of label *I*’s.

and [B+] play against opponents who also have higher incentives than in [A] and [B]. Thus, treatments [A+] and [B+] cannot be directly compared to [E] and [F].

## 5.2 Reasoning about opponents' incentives

In the design of treatments [A+], [B+] and [C+], relative to [A], [B], [C], we increase the payoff for undercutting the opponent for both players in the game. Thus, the shifts in the distributions towards lower numbers observed in Section 4.1 may conflate two distinct effects. The first effect is the possible increase in the cognitive bound of player  $i$ , and the second is the change in  $i$ 's beliefs about  $j$ 's cognitive bound due to the change in  $j$ 's incentives. Both effects would determine an increase in the behavioral  $k_i$ , hence a shift of the distribution towards lower actions.

To illustrate these effects using our model, consider the example in Figure 2.b (p.18). When  $v_i = v_j = v$ , the cognitive bound and behavioral level of play are  $\hat{k}_i = k_i = 3$ . Now, suppose that  $v_i$  is increased up to  $v^*$ , while holding fixed  $v_j = v$ . In this case,  $\hat{k}_i = 6$  and  $\hat{k}_j^i = \bar{k}_j^i = k_j^i = 4$ . Player  $i$ 's response is to play according to  $k_i = 5$ , which is higher than the original level of 3. Suppose next that  $v_j$  is also increased to  $v_j = v^*$ . Then player  $i$ 's cognitive bound becomes binding, and  $k_i$  increases to  $k_i = \hat{k}_i = 6$ . The movement from 3 to 5 is thus due to the increase in  $\hat{k}_i$  alone, induced by an increase in  $v_i$ ; the further change from 5 to 6 instead is determined by  $i$ 's reasoning about the change in his opponent's incentives.<sup>27</sup>

The following treatments, summarized in Table 5.2, are aimed at testing whether subjects in our experiment reason about their opponents' incentives independently of their own. In a sense, the exercise is of a similar spirit to treatment [C], in which subjects play against the number chosen by an opponent in treatment [B]. Similarly, in treatments [K] and [L] agents play the high-payoff game against the number chosen by an opponent in treatments [A] and [C], respectively. Both treatments are administered after the main treatments to all forty subjects from two sessions (one exogenous and one endogenous), and each is repeated three times.

These treatments add a further layer of complexity, since the individual is told in treatment [K] (resp., [L]) that he is playing the high-payoff game against the number chosen by an opponent of the same (other) label himself playing the low payoff game against opponent of the same (other) label. Treatment [L] is especially complex: for player  $i$ , both the payoffs and the label of  $i$ 's opponent *and* of the opponent's opponent are different from  $i$ 's own payoff and label.

By comparing treatments [K] and [L] with treatments [A] and [C] and with treatments [A+] and [C+], we can, in principle, disentangle the two effects mentioned above. The shift from [A] to [K] (and from [C] to [L]), due solely to the increase of each subject's own payoffs and not his opponent's, may be attributed to the increase of subjects' own cognitive bound. It

<sup>27</sup>In general, our model implies that  $\hat{k}_i$  (weakly) increases whenever  $v_i$  is increased. Furthermore, if  $v_j$  is held constant,  $\hat{k}_j^i$  increases only if  $\bar{k}_j^i$  (the intersection between  $c_j^i$  and  $v_j$ ) had been larger than  $\hat{k}_i - 1$  in the first place. In Figure 2.a, for instance, increasing  $v_i$  without changing  $v_j$  does not affect  $\hat{k}_j^i$ , hence  $k_i$ . The opposite is true in Figure 2.b, where  $\hat{k}_j^i$  increases as  $v_i$  is increased, until  $\hat{k}_j^i = \bar{k}_j^i$  (that is, when  $v_i$  is sufficiently high that  $\hat{k}_i \geq 5$ .)



Treatment	Own label	Opponent's label	Own payoffs	Opponent's payoffs	Replacement of opponent's opponent
Mixed payoffs-homogeneous [K]	$I (II)$	$I (II)$	High	Low	No
Mixed payoffs-heterogeneous [L]	$I (II)$	$II (I)$	High	Low	Yes

Table 4: Treatment summary for mixed-payoffs treatments [K] and [L]

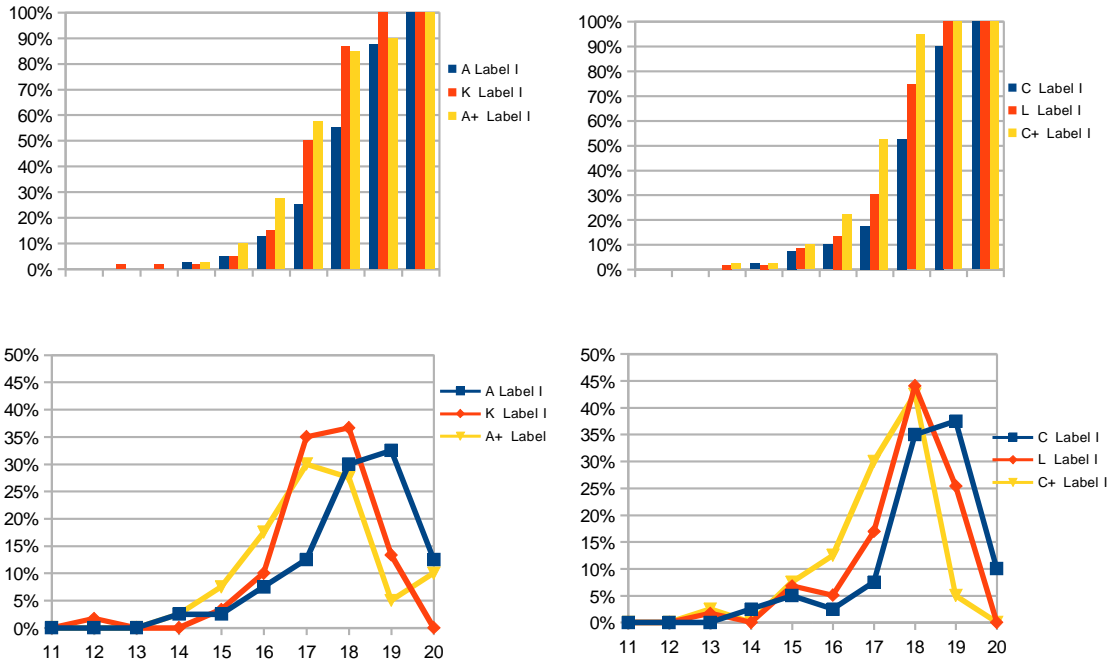


Figure 8: Treatment [K] compared to [A] and [A+]; treatment [L] compared to [C] and [C+] for label  $I$ ; cumulative distributions (top) and frequency distributions (bottom)

should be observed only if the cognitive bound in treatments [A] and [C] had been binding (see footnote 27); the further shift from [K] to [A+] (and from [L] to [C+]) instead can be imputed to the increase in subjects' beliefs about their opponents' behavior due to the increase of their payoffs.

Figures 8 and 9 show the results of these treatments for labels  $I$  and  $II$ , respectively. The results are roughly in line with the predictions of the theory. The empirical distribution of [A] first order stochastically dominates [K] everywhere for label  $I$  other than at 14, and in most of the curve for label  $II$ . Distribution [K] first order stochastically dominates [A+] in the majority of the curve, although this appears more tentative. The results for [L], however, are surprisingly clean: the distribution of [L] lies 'in between' the distributions of [C] and [C+] nearly everywhere for label  $I$  and in most regions for label  $II$ . Furthermore, for label  $I$ , the theory predicts that the increase in  $k_i$  from [C] to [L] should be at most one, which seems roughly confirmed by the small shift in distribution from [C] to [L]. In this case, and consistently with the theory, the movement from [C] to [C+] for label  $I$  is mainly due to the

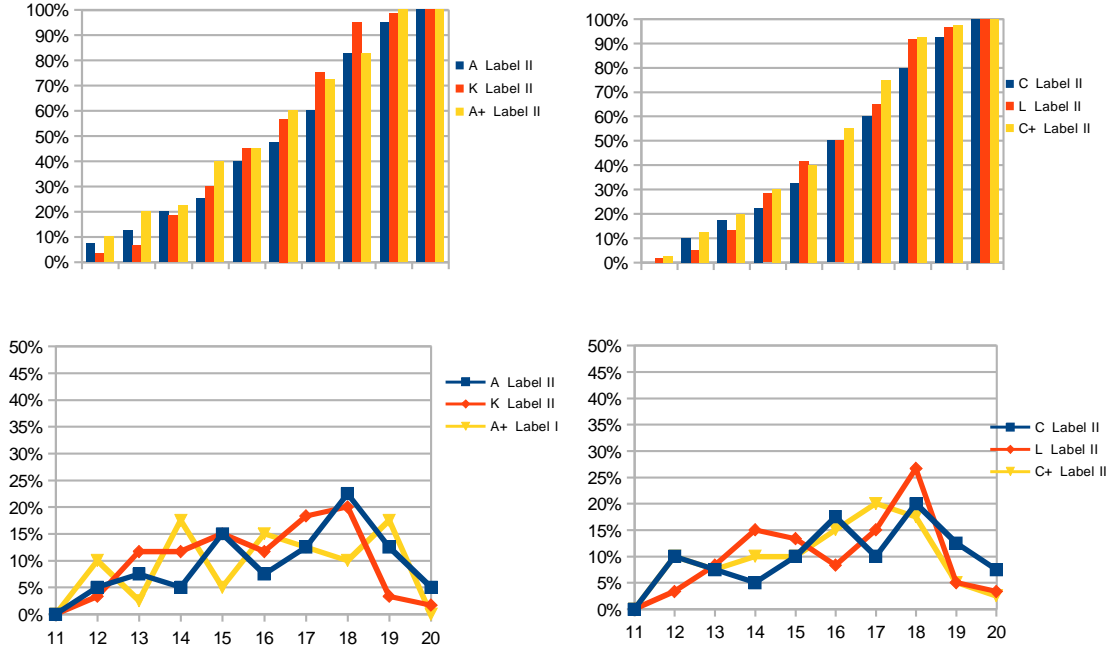


Figure 9: Treatment [K] compared to [A] and [A+]; treatment [L] compared to [C] and [C+] for label II; cumulative distributions (top) and frequency distributions (bottom)

increase in the opponents' payoffs, and not solely to the agent's own incentives. In light of the complexity of these treatments and the difficulty of the instructions, these results appear to be surprisingly good. But these factors suggest caution in interpreting the results.

## 6 Concluding Remarks

In this paper we have extended the level- $k$  approach to strategic reasoning, introducing a model that endogenizes individuals' cognitive bounds as the result of a cost-benefit analysis. Our model further allows a differentiation between players' cognitive bounds and their perception of the cognitive bound of others. Observed behavior is not only a function of game theoretical sophistication but also of incentives and beliefs, illustrating that caution should be exercised in interpreting level of play as purely revealing of cognitive ability.<sup>28</sup> Our framework also solves apparent conceptual difficulties of the level- $k$  approach, such as the possibility that individuals reason about opponents they regard as more sophisticated. The model leads to a rich set of novel predictions which our experiment is designed to test. By jointly accommodating our experimental findings, the model provides a unifying framework of the reasoning process used

<sup>28</sup>In a different setting, it is a well known theme in the Economics of Education literature that incentives may affect standard measures of cognitive abilities. For a recent survey of the vast literature that combines classical economics notions with measurement of cognitive abilities and psychological traits to address the endogeneity problems stemming from the role of incentives, see, for instance, Almlund, Duckworth, Heckman and Kautz (2011).

in initial response games.

In addition to supporting our model, the experiment introduced in this paper plays a more general role. It reveals that individuals change their actions as their incentives and beliefs about the opponents are varied. Moreover, players change their behavior in a systematic way, not captured by existing models of strategic reasoning. Taken together, these findings suggest that there is a fundamental reasoning process behind individual choices. This demonstrates the need for a general framework, and we view our model as being a natural and tractable candidate.

This paper takes a step towards providing a more complete model of procedural rationality under non-equilibrium play; a natural next step is a deeper understanding of the framework. In Alaoui and Penta (2013), we analyze the primitive properties of behavior that lead to the cost-benefit analysis described here, providing an axiomatic foundation to our model. We further characterize the value of reasoning function through additional assumptions, and we apply the model to other games of initial responses. We find that our model can explain the behavior observed in well-known experiments, such as those conducted by Goeree and Holt (2001). We view these results as providing further evidence in support of our theory, and an external validation of our approach.

In closing, we note that by relating individuals' cognitive bounds to their incentives in the game, our theory establishes a link between level- $k$  reasoning and the conventional domain of economics, centered around tradeoffs and incentives. We see this as a desirable feature from a methodological viewpoint, which can further favor the integration of theories of initial responses within the core of economics. Conversely, the application of classical economic concepts to a model of initial responses opens new directions of research from both a theoretical and an empirical viewpoint. For instance, a rigorous identification of the properties of the cost functions would lead to additional predictions, particularly on the magnitudes of the change in play across games. These last questions are outside the scope of this paper, and remain open for future research.

## References

1. Agranov, Marina, Andrew Caplin and Chloe Tergiman. 2012. "Naive Play and the Process of Choice in Guessing Games" *mimeo*.
2. Agranov, Marina, Elizabeth Potamites, Andrew Schotter and Chloe Tergiman. 2012. "Beliefs and Endogenous Cognitive Levels: An Experimental Study" *Games and Economic Behavior*, 75(2): 449-463.
3. Alaoui, Larbi and Antonio Penta. 2013. "Strategic Thinking and the Value of Reasoning: Theory and Applications to Five "Little Treasures" of Game Theory", *mimeo*.
4. Almlund, Mathilde, Angela Lee Duckworth, James Heckman and Tim Kautz. 2011. "Personality Psychology and Economics", *Handbook of the Economics of Education*, Volume 4.
5. Arad, Ayala and Ariel Rubinstein. 2012. "The 11-20 Money Request Game: A Level- $k$  Reasoning Study", *American Economic Review*, 102(7): 3561-3573.
6. Basu, Kaushik, "The Traveler's Dilemma: Paradoxes of Rationality in Game Theory." *American Economic Review Papers and Proceedings*, 84(2): 391-395.
7. Bhatt, Meghana A., and Colin F. Camerer. 2005. "Self-referential Thinking and Equilibrium as States of Mind in Games: fMRI Evidence." *Games and Economic Behavior*, 52(2): 424-459.
8. Bhatt, Meghana A., Terry Lohrenz, Colin F. Camerer, and P. Read Montague. 2010. "Neural Signatures of Strategic Types in a Two-Person Bargaining Game." *Proceedings of the National Academy of Sciences*, 107(46): 19720-19725.
9. Bosch-Domènech, Antoni, Jose García-Montalvo, Rosemarie Nagel, and Albert Satorra. 2002. "One, Two, (Three), Infinity...: Newspaper and Lab Beauty-Contest Experiments." *American Economic Review*, 92(5), 1687-1701.
10. Camerer, Colin F., Teck-Hua Ho, and Juin Kuan Chong. 2004. "A Cognitive Hierarchy Model of Games." *Quarterly Journal of Economics*, 119(3): 861-898.
11. Capra, C. Monica, Jacob K. Goeree, Rosario Gomez, and Charles A. Holt. 1999. "Anomalous Behavior in a Traveler's Dilemma?" *American Economic Review*, 89(3): 678-690.
12. Choi, Syungjoo. 2012. "A Cognitive Hierarchy Model of Learning in Networks." *Review of Economic Design*, 16: 215-250.
13. Coricelli, Giorgio, and Rosemarie Nagel. 2009. "Neural Correlates of Depth of Strategic Reasoning in Medial Prefrontal Cortex." *Proceedings of the National Academy of Sciences*, 106(23): 9163-9168.

14. Costa-Gomes, Miguel A., and Vincent P. Crawford. 2006. "Cognition and Behavior in Two-Person Guessing Games: An Experimental Study." *American Economic Review*, 96(5): 1737-1768.
15. Costa-Gomes, Miguel A., Vincent P. Crawford, and Bruno Broseta. 2001. "Cognition and Behavior in Normal-Form Games: An Experimental Study." *Econometrica*, 69(5): 1193-1235.
16. Crawford, Vincent P. 2003. "Lying for Strategic Advantage: Rational and Boundedly Rational Misrepresentation of Intentions." *American Economic Review*, 93(1): 133-149.
17. 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.
18. Crawford, Vincent P., and Nagore Iriberri. 2007. "Level- $k$  Auctions: Can a Nonequilibrium Model of Strategic Thinking Explain the Winner's Curse and Overbidding in Private-Value Auctions?" *Econometrica*, 75(6): 1721-1770.
19. Davidson, Russell, and Jean-Yves Duclos. 2000. "Statistical inference for stochastic dominance and for the measurement of poverty and inequality." *Econometrica*, 68(6): 1435-1464.
20. Fudenberg, Drew. "Predictive Game Theory." NSF White Papers 2010.
21. Gabaix, Xavier. "Game Theory with Sparsity-Based Bounded Rationality." *mimeo*
22. Goeree, Jacob K., and Charles A. Holt. 2001. "Ten Little Treasures of Game Theory and Ten Intuitive Contradictions." *American Economic Review*, 91(5): 1402-1422.
23. Goeree, Jacob K., and Charles A. Holt. 2004. "A Model of Noisy Introspection." *Games and Economic Behavior*, 46(2): 365-382.
24. Grosskopf, Brit, and Rosemarie Nagel. 2008. "The Two-Person beauty contest." *Games and Economic Behavior*, 62(1): 93-99.
25. Kets, Willemien. 2012. "Bounded Reasoning and Higher-Order Uncertainty." *mimeo*, Kellogg School of Management, Northwestern University.
26. Kneeland, Terri. 2012. "Coordination under Limited Depth of Reasoning." *mimeo*, University of British Columbia.
27. Ho, Teck-Hua, Colin Camerer, and Keith Weigelt. 1998. "Iterated Dominance and Iterated Best Response in Experimental 'p-Beauty Contests'." *American Economic Review*, 88(4): 947-969.

28. Nagel, Rosemarie. 1995. "Unraveling in Guessing Games: An Experimental Study." *American Economic Review*, 85(5): 1313-1326.
29. Palacios-Huerta, Ignacio and Oscar Volij. 2009. "Field Centipedes." *American Economic Review*, 99(4): 1619-1635.
30. Stahl, Dale O., and Paul R. Wilson. 1994. "Experimental Evidence on Players' Models of Other Players." *Journal of Economic Behavior and Organization*, 25(3): 309-327.
31. Stahl, Dale O., and Paul R. Wilson. 1995. "On Players' Models of Other Players: Theory and Experimental Evidence." *Games and Economic Behavior*, 10(1): 218-254.
32. Strzalecki, Tomasz. 2010. "Depth of Reasoning and Higher-Order Beliefs." Harvard Institute of Economic Research Discussion Paper Number 2184.

# Appendix

## A Logistics of the Experiment

The experiment was conducted at the Laboratori d'Economia Experimental (LEEX) at Universitat Pompeu Fabra (UPF), Barcelona. Subjects were students of UPF, recruited using the LEEX system. No subject took part in more than one session. Subjects were paid 3 euros for showing up (students coming from a campus that was farther away received 4 euros instead). Subjects' earnings ranged from 10 to 40 euros, with an average of 15.8.

Each subject went through a sequence of 18 games. Payoffs are expressed in 'tokens', each worth 5 cents. Subjects were paid randomly, once every six iterations. The order of treatments is randomized (see below). Finally, subjects only observed their own overall earnings at the end, and received no information concerning their opponents' results.

Our subjects were divided in 6 sessions of 20 subjects, for a total of 120 subjects. Three sessions were based on the exogenous classification, and each contained 10 students from the field of humanities (humanities, human resources, and translation), and 10 from math and sciences (math, computer science, electrical engineering, biology and economics). Three sessions were based on the endogenous classification, and students were labeled based on their performance on a test of our design. (See Appendix B). In these sessions, half students were labeled as 'high' and half as 'low'.

### A.1 Instructions of the Experiment

We describe next the instructions as worded for a student from math and sciences. The instructions for students from humanities would be obtained replacing these labels everywhere. Similarly, labels high and low would be used for the endogenous classification.

#### A.1.1 Baseline Game and Treatments [A], [B] and [C]

Pick a number between 11 and 20. You will always receive the amount that you announce, in tokens.

In addition:

- if you give the same number as your opponent, you receive an extra 10 tokens.
- if you give a number that's exactly one less than your opponent, you receive an extra 20 tokens.

*Example:*

- If you say 17 and your opponent says 19, then you receive 17 and he receives 19.
- If you say 12 and your opponent says 13, then you receive 32 and he receives 13.
- If you say 16 and you opponent says 16, then you receive 26 and he receives 26.

#### **Treatments [A] and [B]:**

Your opponent is:

- a student from maths and sciences (treatment [A]) / humanities (treatment [B])
- he is given the same rules as you.

#### **Treatment [C]:**

In this case, the number you play against is chosen by:

- a student from humanities facing another student from humanities. In other words, two students from humanities play against each other. You play against the number that one of them has picked.

### A.1.2 Changing Payoffs: Treatments [A+], [B+] and [C+]

You are now playing a high payoff game.

- if you give the same number as your opponent, you receive an extra 10 .

*Example:*

- If you say 17 and your opponent says 19, then you receive 17 and he receives 19.
- If you say 12 and your opponent says 13, then you receive 92 and he receives 13.
- If you say 16 and you opponent says 16, then you receive 26 and he receives 26.

#### Treatments [A+] and [B+]

Your opponent is:

- a student from maths and sciences playing the high payoff game (treatment [A+]) / humanities (treatment [B+])

- he is given the same rules as you.

#### Treatment [C+]

In this case, the number you play against is chosen by:

- a student from humanities playing the high payoff game with another student from humanities. In other words, two students from humanities play the high payoff game with each other (extra 10 if they tie, 80 if exactly one less than opponent). You play against the number that one of them has picked.

### A.1.3 Treatments [D], [E] and [F]

Before playing treatments [D], [E] and F, the subjects were given the following ‘tutorial’:

**Game Theory Tutorial:** According to game theory, if the players are infinitely rational, then the game should be played in the following way. Both players should say 11.

*Explanation:* Suppose the two players are named Ana and Beatriz. If Ana thinks Beatriz plays 20, then Ana would play 19. But then Beatriz knows that Ana would play 19, so she would play 18. Ana realizes this, and so she would play 17.... they both follow this reasoning until both would play 11. Notice that if Beatriz says 11, then the best thing for Ana is to also say 11.

#### Treatment [D]

Your opponent is:

- a student who has also been given the game theory tutorial.

#### Treatment [E]

Your opponent is:

- a student from maths and sciences.
- he has not been given the game theory tutorial.

#### Treatment [F]

Your opponent is:

- a student from humanities.
- he has not been given the game theory tutorial.

### A.1.4 Treatments [K] and [L]

#### Treatment [K]

In this case, the number you play against is chosen by:

- a student from maths and sciences playing the low payoff game with another student from maths and sciences. In other words, two students from maths and sciences play the low payoff game with each



other (extra 10 if they tie, 20 if exactly one less than opponent). You play against the number that one of them has picked.

**Treatment [L]**

In this case, the number you play against is chosen by:

- a student from humanities playing the low payoff game with another student from humanities. In other words, two students from humanities play the low payoff game with each other (extra 10 if they tie, 20 if exactly one less than opponent). You play against the number that one of them has picked.

**A.2 Sequences**

Our 6 groups (3 for the endogenous and 3 for the exogenous classification) went through four different sequences of treatments. Two of the groups in the exogenous treatment followed Sequence 1, and one followed Sequence 2. The three groups of the endogenous classification each took a different sequence: respectively sequence 1, 3 and 4. All the sequences contain our main treatments, [A], [B], [C], [A+], [B+], [C+]. In addition, sequences 2 and 4 contain the [K] and [L] treatments, whereas sequences 1 and 3 conclude with the additional treatments [D], [E] and [F]. The order of the main treatments is different in each sequence, both in terms of changing the beliefs and the payoffs.

- **Sequence 1:**  $A, B, C, B, A, C, A^+, B^+, C^+, B^+, A^+, C^+, D, E, F, D, E, F$
- **Sequence 2:**  $A, B, B, A, C, C, K, L, K, L, K, L, A^+, B^+, B^+, A^+, C^+, C^+$
- **Sequence 3:**  $A^+, B^+, C^+, B^+, A^+, C^+, A, B, C, B, A, C, D, E, F, D, E, F$
- **Sequence 4:**  $B, A, C, B, A, C, K, L, K, L, K, L, B^+, A^+, C^+, B^+, A^+, C^+$

**B The Test for the Endogenous Classification**

The cognitive test takes roughly thirty minutes to complete, and consists of three questions. In the first, subjects are asked to play a variation of the board game Mastermind. In the second question, the subjects are given a typical centipede game of seven rounds, and are asked what an infinitely sophisticated and rational agent would do. In the third game, the subjects are given a lesser known ‘pirates game’, which is a four player game that can be solved by backward induction. Subjects are asked what the outcome of this game would be, if players were ‘infinitely sophisticated and rational’. Each question was given a score, and then a weighted average was taken. Subjects whose score was higher (lower) than the median score were labeled as ‘high’ (‘low’). We report next the instructions of the test, as administered to the students (see the online appendix for the original version in Spanish).

**Instructions of the Test.** This test consists of three questions. You must answer all three, within the time limit stated.

**Question 1:**

In this question, you have to guess four numbers in the correct order. Each number is between 1 and 7. No two numbers are the same. You have nine attempts to guess the four numbers. After each attempt, you will be told the number of correct answers in the correct place, and the number of correct numbers in the wrong place.

*Example:* Suppose that the correct number is: 1 4 6 2.

If you guess : 3 5 4 6, then you will be told that you have 0 correct answers in the correct place and 2 in the wrong place.

If you guess : 3 5 6 4, then you will be told that you have 1 correct answer in the correct place and 1 in the wrong place.

If you guess : 3 4 7 2, then you will be told that you have 2 correct answers in the correct place and 0 in the wrong place.

If you guess : 1 4 6 2, then you will be told that you have 4 correct answers, and you have reached the objective.

Notice that the correct number could not be (for instance) 1 4 4 2, as 4 is repeated twice. You are, however, allowed to guess 1 4 4 2, in any round.

You have a total of 90 second per round: 30 seconds to introduce the numbers and 60 seconds to view the results.

**Question 2:**

Consider the following game. Two people, Antonio and Beatriz, are moving sequentially. The game starts with 1 euro on the table. There at most 6 rounds in this game:

*Round 1)* Antonio is given the choice whether to take this 1 euro, or pass, in which case the game has another round. If he takes the euro, the game ends. He gets 1 euro, Beatriz gets 0 euros. If Antonio passes, they move to round 2.

*Round 2)* 1 more euro is put on the table. Beatriz now decides whether to take 2 euros, or pass. If she takes the 2 euros, the game ends. She receives 2 euros, and Antonio receives 0 euros. If Beatriz passes, they move to round 3.

*Round 3)* 1 more euro is put on the table. Antonio is asked again: he can either take 3 euros and leave 0 to Beatriz, or pass. If Antonio passes, they move to round 4.

*Round 4)* 1 more euro is put on the table. Beatriz can either take 3 euros and leave 1 euro to Antonio, or pass. If Beatriz passes, they move to round 5.

*Round 5)* 1 more euro is put on the table. Antonio can either take 3 euros and leave 2 to Beatriz, or pass. If Antonio passes, they move to round 6.

*Round 6)* Beatriz can either take 4 euros and leaves 2 to Antonio, or she passes, and they both get 3.

Assume Antonio and Beatriz are infinitely sophisticated and rational, and they each want to get as much money as possible. What will be the outcome of the game?

- a) Game stops at Round 1, with payoffs: (Antonio: 1 euro      Beatriz: 0 euros)
- b) Game stops at Round 2, with payoffs: (Antonio: 0 euro      Beatriz: 2 euros)
- c) Game stops at Round 3, with payoffs: (Antonio: 2 euros      Beatriz: 1 euro)
- d) Game stops at Round 4, with payoffs: (Antonio: 1 euro      Beatriz: 3 euros)
- e) Game stops at Round 5, with payoffs: (Antonio: 3 euros      Beatriz: 2 euros)
- f) Game stops at Round 6, with payoffs: (Antonio: 2 euros      Beatriz: 4 euros)
- g) Game stops at Round 6, with payoffs: (Antonio: 3 euros      Beatriz: 3 euros)

You have 8 minutes in total for this question.

**Question 3:**

Four pirates (Antonio, Beatriz, Carla and David) have obtained 10 gold doblónes and have to divide up the loot. Antonio proposes a distribution of the loot. All pirates vote on the proposal. If half the crew or more agree, the loot is divided as proposed by Antonio.

If Antonio fails to obtain support of at least half his crew (including himself), then he will be killed. The pirates start over again with Beatriz as the proposer. If she gets half the crew (including herself) to agree, then the loot is divided as proposed. If not, then she is killed, and Carla then makes the proposal. Finally, if her proposal is not agreed on by half the people left, including herself, then she is killed, and David takes everything.

In other words:

Antonio needs 2 people (including himself) to agree on his proposal, and if not he is killed.

If Antonio is killed, Beatriz needs 2 people (including herself) to agree on her proposal, if not she is killed.

If Beatriz is killed, Carla needs 1 person to agree (including herself) to agree on her proposal, and if not she is killed.

If Carla is killed, David takes everything.

The pirates are infinitely sophisticated and rational, and they each want to get as much money as possible. What is the maximum number of coins Antonio can keep without being killed?

Notice that \*the proposer\* can also vote, and that exactly half the votes is enough for the proposal to pass.

You have 8 minutes in total for this question.

**Scoring.** In the *mastermind* question, subjects were given 100 points if correct, otherwise they received 15 points for each correct answer in the correct place and 5 for each correct answer in the wrong place in their last answer. In the *centipede* game, subjects were given 100 points if they answered that the game would end at round 1, otherwise points were equal to  $\min\{0, (6 - \text{round}) \cdot 15\}$ . In the *pirates* game, subjects obtain 100 if they answer 100, 60 if they answer 10,  $\min\{0, (80 - x) \cdot 10\}$ . The overall score was given by the average of the three.

## C Statistical Tests and Regressions

	Relevant dummy	Classification dummy	Constant	Number of obs.
From A to A+ Label <i>I</i>	-0.50*** (0.17)	0.22 (0.52)	17.21	235
From B to B+ Label <i>I</i>	-0.62*** (0.19)	0.36 (0.38)	17.50	233
From C to C+ Label <i>I</i>	-1.15*** (0.18)	0.34 (0.38)	17.76	236
From A to A+ Label <i>II</i>	-0.64** (0.27)	-0.10 (0.45)	16.91	236
From B to B+ Label <i>II</i>	-0.74*** (0.25)	0.38 (0.47)	16.57	236
From C to C+ Label <i>II</i>	-0.97*** (0.25)	-0.07 (0.45)	16.97	234
From B to A Label <i>I</i>	-0.36* (0.20)	0.34 (0.43)	17.50	236
From B to C Label <i>I</i>	0.25 (0.25)	0.44 (0.37)	17.46	236
From A to C Label <i>I</i>	0.62*** (0.18)	0.36 (0.42)	17.13	236
From A to B Label <i>II</i>	-0.09 (0.26)	0.50 (0.48)	16.6	236
From B to C Label <i>II</i>	0.16 (0.25)	0.36 (0.46)	16.59	235
From A to C Label <i>II</i>	0.07 (0.27)	0.30 (0.44)	16.71	235
From B+ to A+ Label <i>I</i>	-0.26 (0.17)	0.25 (0.46)	16.96	232
From B+ to C+ Label <i>I</i>	-0.29** (0.13)	0.27 (0.45)	16.93	233
From A+ to C+ Label <i>I</i>	0.03 (0.16)	0.21 (0.48)	16.71	236
From A+ to B+ Label <i>II</i>	-0.18 (0.21)	-0.22 (0.52)	16.33	236
From B+ to C+ Label <i>II</i>	-0.07 (0.25)	-0.04 (0.50)	16.05	235
From A+ to C+ Label <i>II</i>	-0.25 (0.21)	-0.47 (0.51)	16.46	235

Table 5: Regressions for Labels *I* and *II*.

	Two-Sample KS D-stat (exact p-value)	Mann-Whitney-Wilcoxon p-values
A vs A+ Label <i>I</i>	0.18 (0.040) **	0.02 **
B vs B+ Label <i>I</i>	0.24 (0.002) ***	0.0004 ***
C vs C+ Label <i>I</i>	0.37 (0.000) ***	0.0000 ***
A vs A+ Label <i>II</i>	0.14 (0.128)	0.019 **
B vs B+ Label <i>II</i>	0.17 (0.048) **	0.0039 ***
C vs C+ Label <i>II</i>	0.25 (0.001) ***	0.0007 ***

Table 6: Equality of Distributions Tests: Changing Payoffs

	Two-Sample KS D-stat (exact p-value)	Mann-Whitney-Wilcoxon p-values
A vs B Label <i>I</i>	0.11 (0.463)	0.21
B vs C Label <i>I</i>	0.05 (0.998)	0.31
A vs C Label <i>I</i>	0.14 (0.163)	0.03 **
A+ vs B+ Label <i>I</i>	0.06 (0.952)	0.47
B+ vs C+ Label <i>I</i>	0.08 (0.738)	0.24
A+ vs C+ Label <i>I</i>	0.07 (0.869)	0.72
A vs B Label <i>II</i>	0.05 (0.986)	0.68
B vs C Label <i>II</i>	0.06 (0.957)	0.76
A vs C Label <i>II</i>	0.04 (1.000)	0.96
A+ vs B+ Label <i>II</i>	0.07 (0.793)	0.42
B+ vs C+ Label <i>II</i>	0.06 (0.938)	0.57
A+ vs C+ Label <i>II</i>	0.12 (0.277)	0.22

Table 7: Equality of Distributions Tests: Changing Opponents

	From A to A+ Label <i>I</i> (t-statistics)	From B to B+ Label <i>I</i> (t-statistics)	From C to C+ Label <i>I</i> (t-statistics)
11	-.0441	-.0967	-.1050
12	-.1344	-.0039	-.1172
13	-.0956	-.0547	-.1401
14	-.0268	-.0298	-.1394
15	-.0449	-.0732	-.2549
16	-.1360	-.1227	-.3846
17	-.2763	-.4009	-.5312
18	-.3625	-.4593	-.6789
19	-.0886	-.1340	-.2222
20			

Table 8: DD test: Changing Payoffs, Label *I*

	From A to A+ Label <i>II</i> (t-statistics)	From B to B+ Label <i>II</i> (t-statistics)	From C to C+ Label <i>II</i> (t-statistics)
11	-.0836	-.0768	-.2209
12	-.0871	-.0635	-.1501
13	-.0412	-.0393	-.0893
14	-.1273	-.0553	-.1858
15	-.0652	-.1313	-.1650
16	-.1338	-.2267	-.1989
17	-.2245	-.2812	-.4879
18	-.2520	-.3309	-.3865
19	-.3090	-.3249	-.2185
20			

Table 9: DD test: Changing payoffs, Label *II*

	Two-Sample KS D-stat (exact p-value)	Mann-Whitney-Wilcoxon p-values
A1 vs A2 Label <i>I</i>	0.06 (0.999)	0.55
B1 vs B2 Label <i>I</i>	0.03 (1.000)	0.90
C1 vs C2 Label <i>I</i>	0.05 (1.000)	0.76
A+1 vs A+2 Label <i>I</i>	0.05 (1.000)	0.65
B+1 vs B+2 Label <i>I</i>	0.11 (0.803)	0.35
C+1 vs C+2 Label <i>I</i>	0.06 (0.999)	0.62
A1 vs A2 Label <i>II</i>	0.10 (0.920)	0.71
B1 vs B2 Label <i>II</i>	0.05 (1.000)	0.96
C1 vs C2 Label <i>II</i>	0.15 (0.412)	0.47
A+1 vs A+2 Label <i>II</i>	0.16 (0.365)	0.13
B+1 vs B+2 Label <i>II</i>	0.084 (0.985)	0.43
C+1 vs C+2 Label <i>II</i>	0.06 (1.000)	0.69

Table 10: Equality of Distributions Tests: Rounds