

Normann, Hans-Theo; Sternberg, Martin

**Working Paper**

## Hybrid collusion: Algorithmic pricing in human-computer laboratory markets

Discussion Papers of the Max Planck Institute for Research on Collective Goods, No. 2021/11

**Provided in Cooperation with:**

Max Planck Institute for Research on Collective Goods

*Suggested Citation:* Normann, Hans-Theo; Sternberg, Martin (2021) : Hybrid collusion: Algorithmic pricing in human-computer laboratory markets, Discussion Papers of the Max Planck Institute for Research on Collective Goods, No. 2021/11, Max Planck Institute for Research on Collective Goods, Bonn,  
<https://hdl.handle.net/21.11116/0000-0008-7A8C-2>

This Version is available at:

<https://hdl.handle.net/10419/245974>

**Standard-Nutzungsbedingungen:**

Die Dokumente auf EconStor dürfen zu eigenen wissenschaftlichen Zwecken und zum Privatgebrauch gespeichert und kopiert werden.

Sie dürfen die Dokumente nicht für öffentliche oder kommerzielle Zwecke vervielfältigen, öffentlich ausstellen, öffentlich zugänglich machen, vertreiben oder anderweitig nutzen.

Sofern die Verfasser die Dokumente unter Open-Content-Lizenzen (insbesondere CC-Lizenzen) zur Verfügung gestellt haben sollten, gelten abweichend von diesen Nutzungsbedingungen die in der dort genannten Lizenz gewährten Nutzungsrechte.

**Terms of use:**

*Documents in EconStor may be saved and copied for your personal and scholarly purposes.*

*You are not to copy documents for public or commercial purposes, to exhibit the documents publicly, to make them publicly available on the internet, or to distribute or otherwise use the documents in public.*

*If the documents have been made available under an Open Content Licence (especially Creative Commons Licences), you may exercise further usage rights as specified in the indicated licence.*

Discussion Papers of the  
Max Planck Institute for  
Research on Collective Goods  
2021/11



**Human-Algorithm Interaction:  
Algorithmic Pricing in Hybrid  
Laboratory Markets**

**Hans-Theo Normann  
Martin Sternberg**

**MAX PLANCK**  
SOCIETY





# **Human-Algorithm Interaction: Algorithmic Pricing in Hybrid Laboratory Markets**

**Hans-Theo Normann / Martin Sternberg**

May 2021

This version: October 2021

# Human-Algorithm Interaction: Algorithmic Pricing in Hybrid Laboratory Markets

**Hans-Theo Normann**

*Düsseldorf Institute for Competition Economics (DICE)  
and Max Planck Institute for Research on Collective Goods, Bonn*

**Martin Sternberg**

*Max Planck Institute for Research on Collective Goods, Bonn*

October 2021

**Abstract:** This paper investigates pricing in laboratory markets when human players interact with an algorithm. We compare the degree of competition when exclusively humans interact to the case of one firm delegating its decisions to an algorithm. We further vary whether participants know about the presence of the algorithm. When one of three firms in a market is an algorithm, we observe significantly higher prices compared to human-only markets. Firms employing an algorithm earn significantly less profit than their rivals. For four-firm markets, we find no significant differences. (Un)certainly about the actual presence of an algorithm does not significantly affect collusion, although humans seem to perceive algorithms as more disruptive.

**JEL classification:** C90, L41.

**Keywords:** algorithms, collusion, human-computer interaction, laboratory experiments

**Contact information.** Normann: Düsseldorf Institute for Competition Economics, Heinrich Heine University Düsseldorf, 40225 Düsseldorf, Germany, normann@dice.hhu.de. Sternberg: Max Planck Institute for Research on Collective Goods, Kurt-Schumacher-Str. 10, 53113 Bonn, Germany, sternberg@coll.mpg.de.

**Acknowledgements.** We are grateful to Christoph Engel, Bernhard Ganglmair, Joe Harrington, Timo Klein, Marcel Schubert, and Frederike Zufall for helpful comments. We also received useful suggestions from seminar participants at ALEA 2021, Bayreuth University, CLEEN Workshop 2021, Competition and Regulation Day Mannheim 2020, DICE Winter School 2020, EARIE 2021, IMPRS Uncertainty Thesis Workshop 2019, MPI Bonn Lab Meeting, and Third Open Antitrust Doctoral Seminar 2019.

# 1 Introduction

Algorithms are increasingly taking over price decisions on behalf of the firms which employ them. Whereas pricing algorithms are also used in more traditional brick-and-mortar retailing, for example in supermarkets<sup>1</sup> or gasoline stations<sup>2</sup>, the strongly growing e-commerce<sup>3</sup> adds to their rapid dissemination. In its “E-commerce Sector Inquiry,” the EU Commission (2017) reports that a majority of online firms track the prices of competitors and two thirds of them use algorithmic software. Hence, algorithms are on the rise.

Algorithms are ideally suited to deal with the wealth of data available online on competitors and customers, but they may also raise antitrust concerns. As a major advantage, algorithms improve internal processes, predict demand and adjust prices. They enable consistent pricing strategies and react immediately to any changes in the market environment (OECD, 2017). While this shows potential benefits for consumer welfare, other studies suggest that algorithms could weaken competition and make supracompetitive prices more likely (Brown and MacKay, 2019; Calvano et al., 2020b, 2021; Klein, 2020). Therefore, the challenges associated with algorithmic pricing, in particular, the risk of tacit collusion are widely discussed (Calvano et al., 2020a; Ezrachi and Stucke, 2016, 2017; Harrington, 2018, 2020; Mehra, 2016; Monopolkommission, 2018; OECD, 2017; Oxera, 2017) and seem to be high on the agenda of competition authorities around the world (British Competition and Markets Authority, 2018, 2021; Bundeskartellamt and Autorité de la Concurrence, 2019; Competition Bureau Canada, 2018).

Our paper contributes to this debate by highlighting the importance of “hybrid” interaction between human players and algorithms. Recently, it has been shown that self-learning algorithms learn to play repeated-game

---

<sup>1</sup>See “Surge Pricing Comes To The Supermarket,” The Guardian, June 4, 2017, available at: <https://bit.ly/3mf9IQp> (last accessed on March 11, 2021).

<sup>2</sup>See Assad et al. (2020) and “Why Do Gas Station Prices Constantly Change? Blame the Algorithms,” Wall Street Journal, May 8, 2017, available at: <https://on.wsj.com/3vRCRo3> (last accessed on March 11, 2021).

<sup>3</sup>In 2020, 72% of internet users in the EU ordered goods or services online. See 2020 Eurostat Community Survey on ICT usage in households and by individuals, available at: <https://bit.ly/3biBEga> (last accessed on March 11, 2021).

strategies that maximize joint profits without explicitly being instructed to do so (Abada and Lambin, 2020; Calvano et al., 2020b, 2021; Klein, 2020).<sup>4</sup> To date, such collusion has only been observed when algorithms compete against algorithms. In many markets, however, algorithms and humans interact with each other (EU Commission, 2017; Wieting and Sapi, 2021). Already, Chen et al. (2016, p. 1348) emphasize that the effects of algorithmic pricing are not yet understood, “especially in heterogeneous markets that include competing algorithmic and non-algorithmic sellers.” More recent studies (Assad et al., 2020; Wieting and Sapi, 2021) likewise suggest that hybrid markets, where humans and algorithms interact, are far from being the exception.

Our main contribution is an analysis of hybrid interaction in a laboratory experiment. We investigate experimental oligopolies when human players interact with an algorithm. We compare the price level when exclusively humans interact to the case when one firm in the market delegates its decisions to an algorithm. Our research question is whether the presence of an algorithm leads to an increase in prices.<sup>5</sup>

To explore the role of human beliefs about algorithms, we further vary (in a non-deceptive manner) whether or not participants know about the presence of the algorithm. Do participants behave differently when they are aware they are facing an algorithm? This may indeed be the case: Studies on “algorithm aversion” show that people avoid algorithmic advice even though the algorithm is superior to humans (Dietvorst and Bharti, 2020; Dietvorst et al., 2016, 2015). Furthermore, Farjam and Kirchkamp (2018)

---

<sup>4</sup>See Klein (2020) for a comprehensive literature survey on self-learning algorithms, especially those that reply on Q-learning. See also Waltman and Kaymak (2008) for an early dynamic programming approach.

<sup>5</sup>Surprisingly few laboratory experiments have studied cooperation when one or more players are computerized, and none compares human-computer to all-human interaction. Roth and Murnighan (1978) and Murnighan and Roth (1983) analyze two-player prisoner’s dilemmas when subjects know they face a programmed opponent (the papers are well known to be the first to study “infinitely” repeated games in the lab by imposing a random move that determines the end of a supergame). Duffy and Xie (2016) have single humans play against  $n - 1$  computerized grim trigger players, and the authors vary  $n$ . Recently, Duffy et al. (2021) let participants play two-player prisoner’s dilemma supergames against a grim-trigger robot. In Duffy and Xie (2016) and Duffy et al. (2021), subjects know the strategy the robot plays. This is not the case in Roth and Murnighan (1978) and Murnighan and Roth (1983). Again, all four studies have in common that they do not have comparison treatments when the opponents are human.

show in a laboratory asset market that humans trade differently if they expect algorithmic traders. As a possible explanation, they suggest that human traders perceive the algorithmic traders as behaving more rationally. De Melo et al. (2015) find that people tend to make different decisions depending on whether they are facing a human or a computer algorithm.<sup>6</sup> Thus, it seems warranted to test whether expectations about algorithms influence the behavior of participants.

To analyze these research questions, we opted for a rather simple and transparent experimental design. In three- and four-firm markets,<sup>7</sup> participants have two actions (high price, low price) available, so they play an  $n$ -player prisoner’s dilemma. As mentioned, one of the human participants may be replaced by an algorithm. The algorithm we use is a multiplayer generalization of tit-for-tat (Axelrod, 1984; Hilbe et al., 2015). It begins by cooperating, but subsequently adapts to the level of cooperation in the market. Tit-for-tat is cooperative, so when matched with other cooperative strategies, it achieves collusive payoffs. It is also forgiving in that it can return to cooperation after an accidental deviation. It furthermore avoids the exploitation by defectors.

We chose this specific algorithm in order to give cooperation a reasonable good chance. Our algorithm is comparable to the relatively simple programs used in online markets where static algorithms that follow a manageable number of simple rules appear to be common (British Competition and Markets Authority, 2018; Monopolkommission, 2018). Wieting and Sapi (2021) analyze the e-commerce platform *Bol.com* (the largest online marketplace in Belgium and the Netherlands) and identify pricing software that was foremost “relatively unsophisticated” and “consist of a finite set of

---

<sup>6</sup>For related findings, see Dijkstra et al. (1998), Weibel et al. (2008), Krach et al. (2008), Lee (2018) and Rilling et al. (2004).

<sup>7</sup>Experiments with three and four firms seem promising when it comes to identifying collusive effects in that duopolies can be collusive, whereas markets with four or more firms are usually not, see Engel (2015), Fonseca and Normann (2012), Huck et al. (2004), Potters and Suetens (2013). The evidence on cooperation in three-player groups is somewhat inconclusive (and hence a good starting point for us). While Horstmann et al. (2018) do find some collusion in  $n = 3$  oligopolies with differentiated goods, Freitag et al. (2020) do not find any supracompetitive outcomes in a multimarket context benign to collusion. Already Marwell and Schmitt (1972) reported that three-person prisoner’s dilemmas are substantially less cooperative than two-player experiments. Roux and Thöni (2015) demonstrate that larger oligopolies become collusive only when targeted punishments are available.

if-then statements.” Key to our research, British Competition and Markets Authority (2018), Crandall et al. (2018) and Musolff (2021) assume this lack of sophistication may actually increase the risk of supracompetitive prices, as they suspect that particularly simple algorithms lead to collusion. In line with this, Wieting and Sapi (2021) conclude that “[a] secret to successful collusion may lie in managers’ ability to commit to simple strategies.” With self-learning algorithms, our algorithm has in common that it is memory one.<sup>8</sup> Dal Bó and Fréchette (2019) and Romero and Rosokha (2019) recently found that the strategies of human subjects in lab experiments are often memory-one. At the same time, we emphasize that our algorithm is not ferociously committed to cooperate and thus seems suitable to meaningfully study human-algorithm interaction.

Our findings are as follows. Three-firm markets involving an algorithmic player are significantly more collusive than human-only triopolies. While this higher level of collusion raises profits for all firms in the industry, it turns out that those firms that employ the algorithm earn significantly less profit than their rivals. Thus, one firm must be willing to accept the set-up costs for the collusion. For four-firm markets, we find no significant differences. Knowing or not knowing about the presence does not affect competition significantly. Interestingly, however, humans seem to link cooperation to human behavior and not an algorithm.

---

<sup>8</sup>In their online appendix, (Calvano et al., 2020b) briefly report on memory-two algorithms. As the state space disproportionately increases with a two-period memory, these algorithms perform less collusive and have substantially longer punishment phases.



## 2 Experimental Design

The stage game underlying the experiment is an  $n$ -player prisoner’s dilemma,  $n \in \{3, 4\}$ , framed as a market interaction. Players choose a high price or a low price, so the action set for all players is  $\{p_{high}, p_{low}\}$ . Payoffs for  $n = 3$  are as in Table 1 (which is similar to the one used in the experiment). These payoffs are derived from a Bertrand oligopoly model with inelastic demand and constant marginal costs of production.<sup>9</sup> For  $n = 3$  players and actions  $p_{high} = 100$  and  $p_{low} = 60$ , the payoffs in Table 1 result. The payoff table for the four-player treatments can be found below in Section 5.3.

		Other firms’ prices		
		$p_{high}, p_{high}$	$p_{high}, p_{low}$	$p_{low}, p_{low}$
Own price	$p_{high}$	800	0	0
	$p_{low}$	1,440	720	480

Table 1: Payoff table ( $n = 3$  treatments).

We compare six different treatments. We run four treatments with  $n = 3$  players, and we conduct two further treatments with  $n = 4$  players. We vary treatments with and without algorithms and treatments with and without information on the presence of the algorithm. See Table 2.

In all experiments, groups of  $n \in \{3, 4\}$  participants constitute one market. In the treatments labeled “Human\_,” there are  $n$  human players. In the treatments labeled “Algorithm\_,” there are  $n - 1$  human players and one algorithm. In treatments involving an algorithm, the computer decides on behalf of one human; the  $n^{th}$  human is an experimental subject, but he or she is inactive and merely obtains the payoff earned by the algorithm.

The second treatment dimension indicates whether the participants

---

<sup>9</sup>Suppose there are  $m = 24$  consumers who demand one unit of the good up to a reservation price of 100. Each player can supply all consumers at production costs of zero. The player that charges the lowest price serves all consumers; if several players charge the lowest price, they split the profit equally.

$n$ humans	$n - 1$ humans 1 algorithm
Human_Uncertain_3	Algorithm_Uncertain_3
Human_Certain_3	Algorithm_Certain_3
Human_Certain_4	Algorithm_Certain_4

Table 2: Treatment design.

know the composition of the market.<sup>10</sup> These treatments were done with  $n = 3$  players only. In the treatments labeled “\_Certain,” participants know from the instructions whether or not an algorithm is present. In the “\_Uncertain” treatments, the participants do not know if they are part of the Human\_Uncertain or the Algorithm\_Uncertain treatment, so they do not know whether an algorithm is present. They are merely told that, with a probability of 50%, one of the three subjects’ decisions is taken by an algorithm. We conducted the same number of sessions in both treatments. Thus, consistent with the instructions, there was a 50% chance that the participants were in the Algorithm\_Uncertain\_3 treatment.

The algorithm is programmed to play *proportional tit-for-tat*, or *pTFT* (Hilbe et al., 2015). It is an  $n$ -player generalization of tit-for-tat (Axelrod, 1984): Let  $t$  be the index for time. The algorithm begins by cooperating in the first period ( $t = 0$ ) and later cooperates proportionally to the number of cooperators in the previous period. Accordingly, *pTFT* chooses the high price with the following probabilities

$$prob.(p = p_{high}) = \begin{cases} 1 & \text{if } t = 0 \\ \frac{j}{n-1} & \text{if } t > 0 \end{cases}$$

where  $n$  is the number of players including the algorithm player and  $j \in \{0, 1, 2, \dots, n - 1\}$  is the number of rival players who chose  $p_{high}$  in the previous period. Subjects are not told how the algorithm is programmed. Nor are they told the algorithm’s purpose.

<sup>10</sup>Regarding this point, our design is similar to the one in Farjam and Kirchkamp (2018).

The treatments are implemented as repeated games, and all treatments have three supergames.<sup>11</sup> The subjects stay in the same market throughout the periods of the supergames. When a new supergame begins, subjects are randomly assigned to a new market. In other words, we have fixed matching within supergames and random matching across supergames. Each supergame lasts at least 20 periods. From the 20th period onward, a random rule with a continuation probability of 7/10 determines whether play continues. The number of periods in all three rounds was determined ex ante and is the same in all sessions (24, 20 and 21 periods for  $n = 3$ ; and 22, 25, 21, for  $n = 4$ ). Subjects knew they would play three supergames from the instructions and they also knew the termination probability.

### 3 Hypotheses

An algorithm may affect human behavior via (at least) two channels. We call these the *belief channel* and the *action channel*. The two channels affect behavior differently and give rise to different hypotheses.<sup>12</sup>

We begin with actions. At least in the long run, human subjects will probably be influenced by the algorithm’s actual price-setting behavior and its responses, including the punishments it triggers, and so on. In other words, the algorithm’s actions will matter.

Our *pTFT* algorithm is more collusive than the average human and this should have a positive effect on the proportion of  $p_{high}$  choices in a market. A prominent recent literature (Dal Bó and Fréchette, 2011, 2018; Bigoni et al., 2015; Fudenberg et al., 2012) investigates in two-player PDs repeated-game *strategies* and finds that the three strategies *always defect* (*AD*), *grim trigger* (*GT*), and *tit-for-tat* (*TFT*) account for most of the data. The *pTFT* algorithm always cooperates in the first period and rewards cooperation in the following periods. So the algorithm is somewhat

---

<sup>11</sup>See Honhon and Hyndman (2020) for an analysis of how matching schemes and reputation mechanisms affect cooperation in the repeated prisoner’s dilemma.

<sup>12</sup>A third channel could be altered other-regarding preferences: Participants may feel inclined to defect when playing with an algorithm, but not with a human participant, especially if the money earned by the algorithm is kept by the experimenter. Our Algorithm\_ treatments, however, involved three human participants; the profit earned by the algorithm was paid out to a (passive) human participant. Therefore, altered other-regarding preferences should not play a role.

forgiving and willing to resume to cooperation, provided other players do so. This is in contrast to *GT* (and of course *AD*).<sup>13</sup> Thus, compared to a human who plays *AD*, *GT* or *TFT* with some probability, the algorithm is more cooperative. Modelling a players decision, we show in the Appendix A.1 that cooperation in the presence of an *pTFT*-algorithm can be a subgame-perfect Nash equilibrium. Furthermore, if we take strategic risk into account, we demonstrate that the minimum discount factor required for cooperation is higher for three *GT* players compared to two *GT* players and one *pTFT* (Appendix A.1). Altogether, we hypothesize the *pTFT* algorithm should increase cooperation.

We turn to beliefs. Human subjects may expect the algorithm to play differently than other humans. Responding to this belief, humans adjust their behavior accordingly.<sup>14</sup> But in which direction will the belief be affected?

We hypothesize humans to be skeptical about the play of an algorithmic competitor, so they expect less cooperation when an algorithm is present. News about algorithms beating humans at Chess or Go demonstrate the power of machines in zero-sum games. This may suggest that humans also lose against the algorithm in the market domain – that is, firms run by humans earn less profit. Trust is an important part of successful collusion, but the literature on algorithm aversion (Dietvorst and Bharti, 2020; Dietvorst et al., 2015) suggests that humans trust algorithms less than other humans. Along these lines, Farjam and Kirchkamp (2018) find that algorithms are perceived as “more rational.” This could correspond to skeptical expectations. From the subjects’ perspective, “rationality” could imply that the algorithm will attempt to exploit human participants to gain higher profits through competitive behavior. Reports on competitive algorithmic price wars<sup>15</sup> and the fact that online shopping – often asso-

---

<sup>13</sup>Strategies that are more lenient than *TFT* include tit-for-two-tats and tit-for-three-tats. They defect only after the other player has defected two/three times. And perfect tit-for-tat (or win-stay-lose-shift) is a strategy that, unlike *TFT*, even actively returns to cooperation after an all-defect outcome.

<sup>14</sup>There is ample evidence that human subjects respond to beliefs about the action of others. In the prisoner’s dilemma, there are two motives for defection (Ahn et al., 2001; Blanco et al., 2014; Charness et al., 2016). Subjects fear being exploited by others, but some may greedily also want to exploit others themselves.

<sup>15</sup>For example, CBNC reports on undercutting competition between Wal-Mart and Amazon through algorithmic pricing: Sarah Whitten, “Wal-Mart Scammed Into Selling

ciated with algorithmic pricing – is considered a low-price alternative to brick-and-mortar purchases<sup>16</sup> could give rise to the notion that algorithms are particularly competitive. Furthermore, recall that participants are unaware of the algorithm’s strategy, so there is little to suggest that subjects think it is pursuing long-run joint-profit maximization. We thus expect that subjects perceive algorithms as less cooperative than humans and this should, *ceteris paribus*, yield lower cooperation rates.

We now put these conjectures together and state our hypotheses. We begin with the influence of the algorithm’s action. The `_Uncertain` treatments are identical in terms of the instructions and the possibility of an algorithm being present, so the beliefs cannot matter. Here, however, the algorithm’s actual play may have an impact. We hypothesize:

**Hypothesis H<sub>1</sub>:** Cooperation rates in `Algorithm_Uncertain_3` are higher than those in `Human_Uncertain_3`.

For the `_Certain` treatments, the actual play of the algorithm, on the one hand, and skeptical beliefs, on the other, imply ambiguous effects of algorithms. We hypothesize that the use of algorithms will have a positive overall impact on collusion because we provide ample evidence of learning (three relatively long repeated games). Given these learning opportunities, subjects may update their beliefs and adjust them according to the more cooperative behavior of the algorithm. We hypothesize:

**Hypothesis H<sub>2</sub>:** Cooperation rates in `Algorithm_Certain_3` are higher than those in `Human_Certain_3`.

Next, we consider the two `Algorithm_` treatments. The algorithm’s

---

PlayStation 4 for \$90” CNBC November 18, 2014, available at: <https://cnb.cx/3BiNRfm> (last accessed on March 11, 2021); Also, the consultancy Simon-Kucher & Partners reports, in its 2019 Global Pricing Study, that 57% of the companies report they are currently involved in a price war. Available at: <https://bit.ly/3mjWt15> (last accessed on March 11, 2021).

<sup>16</sup>Prices play an important role in online shopping. Degeratu et al. (2000) find that online promotions are stronger signals for price discounts than offline promotions and the price sensitivity of consumers is higher online. A representative survey of German consumers has also shown that 52% of them are convinced that it is cheaper to buy products online, available at: <https://bit.ly/2Zu7PGv> (last accessed on March 11, 2021).

actions are the same here, but in *Algorithm\_Certain*, subjects know for sure they are facing an algorithm, whereas in *Algorithm\_Uncertain*, they might still be competing with a human. Based on our presumption of skeptical beliefs, we hypothesize:

**Hypothesis H<sub>3</sub>:** Cooperation rates in *Algorithm\_Uncertain\_3* are higher than those in *Algorithm\_Certain\_3*.

For the *Human\_* treatments, it is the other way round. The third player is controlled by a human either way, but in *Human\_Uncertain*, participants expect to meet an algorithm with 50% likelihood. So participants should be more optimistic in *Human\_Certain*. We obtain:

**Hypothesis H<sub>4</sub>:** Cooperation rates in *Human\_Certain\_3* are higher than those in *Human\_Uncertain\_3*.

Taking hypotheses H<sub>1</sub> to H<sub>4</sub> together, we obtain an unambiguous and testable ranking for the cooperativeness of our treatments. We should observe that the levels of cooperation satisfy

$$Algorithm_U_3 > Algorithm_C_3 > Human_C_3 > Human_U_3 \quad (1)$$

Finally, we look at markets with four competitors. Here, we do not maintain a directed hypothesis regarding *Algorithm\_Certain\_4* and *Human\_Certain\_4*. On the one hand, *Algorithm\_Certain\_4* should be more collusive than *Human\_Certain\_4* due to the presence of the algorithm. On the other hand, markets with  $n = 4$  firms may already be too competitive for significant collusion to occur at all. Collusive outcomes in market experiments are correlated with the number of firms, and prior experiments indicate that four competitors rarely reach collusive conduct (Engel, 2015; Fonseca and Normann, 2012; Horstmann et al., 2018; Huck et al., 2004; Potters and Suetens, 2013). For the sake of completeness, we state the Null here:

**Hypothesis H<sub>5</sub>:** (Null) Cooperation rates in *Algorithm\_Certain\_4* do not differ from those in *Human\_Certain\_4*.

## 4 Procedures

Subjects were recruited from pools of subjects who had previously volunteered to participate in lab experiments. The experiments involved 429 participants in total. None of the subjects participated in more than one session. We had 24 sessions in total, four for each of the six treatments. Due to different market sizes (three and four players) and participants not showing up, session sizes varied between 12 and 30 participants. The experimental sessions were conducted at labs in Düsseldorf and MPI Bonn between August 2019 and October 2020. No sessions were conducted between early March and mid-July 2020, due to the pandemic. Sessions from mid-July 2020 on were conducted under common hygiene rules. See A.5 in the appendix for session details.

Upon arrival at the laboratory, subjects were randomly assigned to a cubicle, using tokens with the cubicle numbers. After a sufficient number of participants had arrived, the experiment started and participants received a hard copy of the instructions in German. While reading the instructions, subjects were allowed to ask questions privately in their cubicles. Afterwards, control questions made sure everyone had understood the task.

The decision-making parts were conducted as follows. We programmed the experiment in z-Tree (Fischbacher, 2007). In each period, the subjects had to decide by clicking a button whether they wanted to set  $p_{high}$  or  $p_{low}$ . After everyone had decided, an information screen displayed the choices of all three firms in the market and informed subjects about their payoff. At the end of a supergame, the individual overall payoff for that supergame was displayed and the subjects were informed that they would now be assigned to a new market, unless it was the last supergame.

We used an Experimental Currency Unit, where 1,000 ECU corresponded to 1 Euro. One of the three supergames was randomly chosen for payout. At the end of the third supergame, the subjects were informed about the supergame selected for payout and their total earnings.

In the `_Uncertain` treatments, we further asked participants whether they thought an algorithm was present in the experiment. This was done at the end after the last period of the last supergame. Subjects had to

enter a number between zero and 100, expressing how confident they were that an algorithm was in the market. They were paid up to 2 euros for a correct guess: Given a guess  $x \in \{0, 1, 2, \dots, 100\}$  that an algorithm was present, the payoff was  $2x/100$  if an algorithm was actually present and  $2 - 2x/100$  if not. (Participants for whom the algorithm decided were paid 1 euro flat instead.)

The sessions lasted for about 60 minutes. The average payment was 16.8 euro, including a show-up fee.

## 5 Results

### 5.1 Overview

Figure 1 shows how cooperation rates<sup>17</sup> in the different treatments develop over time and supergames.<sup>18</sup> Generally, cooperation increases across supergames: In supergame 1, cooperation rates vary roughly between zero and less than 30%, whereas in supergame 3 they vary between 20 and more than 50%. It appears participants learn to collude tacitly with repetitions of the supergame, confirming the results of Bigoni et al. (2015), Dal Bó and Fréchette (2011, 2018), and Fudenberg et al. (2012).

A closer look reveals that cooperation rates improve for all treatments in supergame 2, but when comparing supergames 2 and 3, only the treatments involving an algorithm increase substantially.<sup>19</sup> This is tentative evidence that the algorithm has a collusive impact.

Complementing Figure 1, Table 3 shows the cooperation rates averaged across periods 6 to 19. We note that the Algorithm\_ treatments have higher averages than their Human\_ counterparts for all treatments in (almost) all supergames.<sup>20</sup> Taking all supergames into account, the highest

---

<sup>17</sup>The cooperation rate is defined as the number of  $p_{high}$  choices divided by the total number of choices, given a treatment or period of play.

<sup>18</sup>To exclude the pronounced restart and endgame effects we observe, we focus on periods 6 to 19. The same graph including all periods can be found in A.6 of the appendix.

<sup>19</sup>There is a very minor increase of cooperation in Human\_Certain\_3 by 0.2 percentage points when comparing supergames 2 and 3. See Table 3.

<sup>20</sup>The exception is that Human\_Certain\_4 is slightly more cooperative than Algorithm\_Certain\_4 in supergame 2.



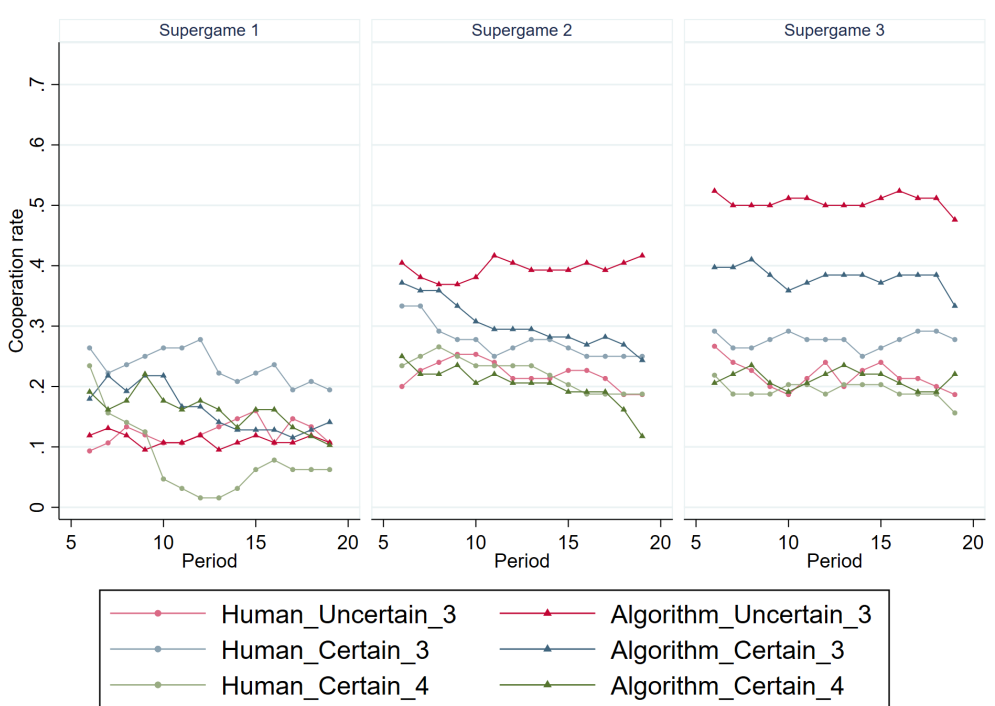


Figure 1: Cooperation rates over time (periods 6 to 19).

cooperation rate is observed in Algorithm.Uncertain.3 (0.337), followed by Algorithm.Certain.3 (0.282) which, in turn, exhibits more cooperation than Human.Certain.3 (0.262). We find higher cooperation in Human.Certain.3 than in Human.Uncertain.3 (0.187). This is exactly the ranking of treatments we hypothesize in (1). Also, in Algorithm.4 (0.191) and Human.4 (0.165), the cooperation rates are lower than in Algorithm.3 and Human.3. This order does not change if we include all periods or focus only on the decisions of human subjects (that is, if we exclude the algorithms' decisions). See A.6 and A.7 of the appendix for details.

How successful are the firms in actually establishing the collusive outcome? Figure 2 shows the percentages of three outcomes for the six treatments in the last supergame: “successful collusion” indicates (tacit) cooperation – all firms choose  $p_{high}$ ; “competition” means that all firms charge  $p_{low}$ ; and “failed collusion” occurs when at least one firm chooses  $p_{low}$  and at least one firm tried to collude – this is miscoordination. Again, it becomes clear that successful coordination on the high price occurs more often in Algorithm.Uncertain.3 and Algorithm.Certain.3. The two extremes

	<i>Supergame 1</i>	<i>Supergame 2</i>	<i>Supergame 3</i>	<i>All</i>
Human_U_3	0.123 (0.328)	0.221 (0.415)	0.218 (0.413)	0.187 (0.390)
Algorithm_U_3	0.111 (0.315)	0.395 (0.489)	0.506 (0.500)	0.337 (0.473)
Human_C_3	0.233 (0.423)	0.275 (0.447)	0.277 (0.448)	0.262 (0.440)
Algorithm_C_3	0.162 (0.369)	0.303 (0.460)	0.381 (0.486)	0.282 (0.450)
Human_C_4	0.0804 (0.272)	0.222 (0.416)	0.193 (0.395)	0.165 (0.371)
Algorithm_C_4	0.160 (0.366)	0.202 (0.401)	0.212 (0.409)	0.191 (0.393)

Standard deviations in parentheses.

Table 3: Average cooperation rates (periods 6 to 19).

are Algorithm\_Uncertain\_3 with a roughly 50% rate of successful collusion, whereas Human\_Certain\_4 involved almost 80% competition. The share of outcomes with miscoordination (failed collusion) is remarkably small in all treatments, meaning that subjects quickly coordinate on either the cooperative or the competitive outcome. This is also apparent from the quick drop in cooperation in the first five periods (see in A.6 of the appendix).

We now systematically test our hypotheses and make statistically reliable statements about treatment effects. We begin with the three-firm markets (5.2). The results from markets with four firms follow in Section 5.3. Throughout, we take the possible dependence of observations into account using bootstrapping standard errors at the session level. See Cameron et al. (2008).

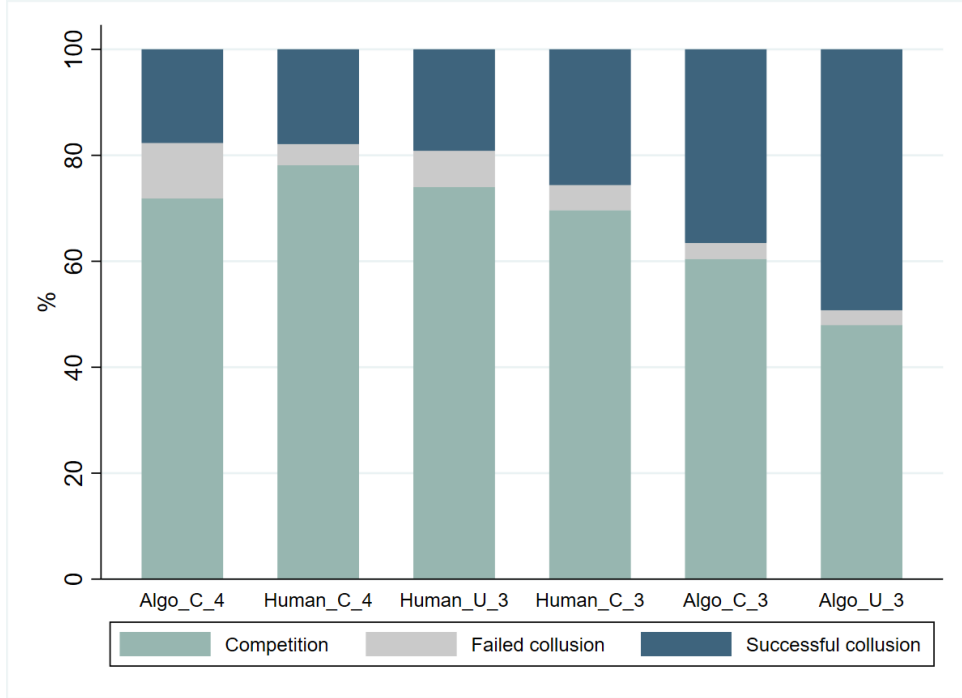


Figure 2: Collusive and competitive outcomes (supergame 3, periods 6 to 19).

## 5.2 Three-firm treatments

### 5.2.1 Treatment differences

Table 4 shows the results of a linear probability model and highlights the main treatment effects we observe in the markets with three firms. Our dependent variable is whether or not a firm (participant or algorithm) cooperates in a given period. We include as explanatory variables dummies for the `Algorithm_` treatments, the `_Certain` treatments, the interaction of the two, and for the initial and terminal periods of play. We report the results separately for the three supergames and jointly for all supergames where we add a cardinal variable for supergame. For the regression including all supergames, the constant reflects supergame 1.

The impact of the algorithm is positive and statistically significant from supergame 2 onward. When analyzing them jointly, the two `Algorithm_` treatments cooperate better than the two `Human_` treatments. The impact of `_Certain` is positive, small, and insignificant. When we add the interaction  $algorithm \times certain$ , the coefficient  $algorithm$  becomes stronger

and remains significant. This indicates that the cooperation rates in Algorithm\_Uncertain\_3 and Human\_Uncertain\_3 differ significantly (Hypothesis H<sub>1</sub>). Comparing the two \_Certain\_3 treatments separately, we find no significant effect of the algorithm, so no support for Hypothesis H<sub>2</sub>.<sup>21</sup>

**Result 1.** *In the  $n = 3$  markets, the Algorithm\_ treatments jointly exhibit significantly higher prices than the Human\_ treatments. Cooperation rates are significantly higher in Algorithm\_Uncertain\_3 compared to Human\_Uncertain\_3. We find no statistically significant effects when comparing Algorithm\_Certain\_3 and Human\_Certain\_3.*

We hypothesize that human subjects play more competitively if they knowingly face or expect to face an algorithmic opponent. But neither the comparison of Algorithm\_Uncertain\_3 and Algorithm\_Certain\_3 (H<sub>3</sub>), nor of Human\_Uncertain\_3 and Human\_Certain\_3 (H<sub>4</sub>) shows significant effects. This suggests that expectations do not play a major role in these regressions.

**Result 2.** *We find no statistically significant effects between the \_Uncertain\_3 and the \_Certain\_3 treatments.*

One interpretation of Result 2 is that expectations do not matter much when subjects gain experience. Below, we report on period-one data, but even in the first period we cannot find much statistical support regarding differences between \_Uncertain\_3 and \_Certain\_3. It appears that even in the very first period of play (first period of the first supergame), when subjects are inexperienced, beliefs do not have a big impact.

---

<sup>21</sup>Across all supergames, the effect of the algorithm in the \_Certain\_3 treatments is statistically insignificant ( $p > 0.1$ ).

	<i>Supergame 1</i>		<i>Supergame 2</i>		<i>Supergame 3</i>		<i>All</i>	
algorithm	-0.0225 (0.0528)	0.00490 (0.0566)	0.108** (0.0526)	0.172*** (0.0654)	0.185** (0.0813)	0.255** (0.110)	0.0847 (0.0516)	0.137** (0.0556)
certain	0.0706 (0.0542)	0.100 (0.0966)	0.000116 (0.0515)	0.0691 (0.0839)	-0.0292 (0.0833)	0.0466 (0.112)	0.0166 (0.0523)	0.0733 (0.0818)
algorithm × certain		-0.0564 (0.117)		-0.132 (0.0991)		-0.145 (0.170)		-0.108 (0.105)
periods 1 to 5	0.0867*** (0.0234)	0.0867*** (0.0235)	0.129*** (0.0177)	0.129*** (0.0173)	0.116*** (0.0199)	0.116*** (0.0198)	0.110*** (0.0133)	0.110*** (0.0132)
periods 20 to 25	-0.0524*** (0.0159)	-0.0524*** (0.0159)	-0.0522*** (0.0160)	-0.0522*** (0.0159)	-0.129*** (0.0292)	-0.129*** (0.0289)	-0.0821*** (0.0108)	-0.0821*** (0.0107)
supergame							0.0968*** (0.0216)	0.0968*** (0.0210)
Constant	0.133*** (0.0405)	0.118*** (0.0372)	0.245*** (0.0518)	0.211*** (0.0561)	0.268*** (0.0625)	0.231*** (0.0531)	0.120*** (0.0401)	0.0923*** (0.0309)
Obs.	7,416	7,416	6,180	6,180	6,489	6,489	20,085	20,085
$R^2$	0.025	0.027	0.029	0.033	0.057	0.063	0.062	0.066

Standard errors in parentheses

\*\*\* p<0.01, \*\* p<0.05, \* p<0.1

Table 4: Treatment effects,  $n = 3$  variants, all periods, linear probability model.

Support for our beliefs hypothesis can nevertheless be detected. The mean cooperation rates correspond to our hypothesis. For the Algorithm\_ treatments, we find (insignificantly) more collusion in \_Uncertain than in \_Certain in supergames two and three and all supergames. For the Human\_ variants, subjects were more collusive in \_Certain than in \_Uncertain in all supergames, as predicted. Recall that we can rank our  $n = 3$  treatments according to our hypotheses, see (1). Creating a cardinal rank variable (with 1 = Algorithm\_Uncertain and 4 = Human\_Uncertain) in our dataset allows for an ordered alternative hypothesis of the multiple independent samples jointly.<sup>22</sup> We see that this ranking variable has a significantly negative effect on the choice in the third supergame (linear probability model,  $p < 0.05$ , see Appendix A.2). The effect is also significant for the second supergame and over all supergames (both  $p < 0.05$ ).

**Result 3.** *Consistent with our hypotheses, the variable ranking for the order of competitiveness of the treatments has a significant negative effect on the cooperation rate.*

Another piece of evidence in favor of our hypothesis on expectations comes from the incentivized guess in the \_Uncertain treatments. This is what we analyze in detail next.

### 5.2.2 Beliefs about the presence of an algorithm

In the two \_Uncertain treatments, we asked participants in an incentivized manner at the end of the experiment about their beliefs of whether one of the firms was equipped with an algorithm. Subjects had to state a probability (a number between zero and 100) that an “algorithm was present in the experiment.”

It turns out that subjects in Human\_Uncertain maintain an average belief of 58.05%, whereas those in Algorithm\_Uncertain have a belief of 45.87%. Table 5 shows the results of a linear probability regression with data from Algorithm\_Uncertain\_3 and Human\_Uncertain\_3 and *algorithm* as an explanatory variable. The variable *algorithm* is significant at  $p <$

---

<sup>22</sup>Similar to the non-parametric Jonckheere-Terpstra test, which is likewise highly significant.

	<i>Guess</i>	
algorithm	-12.19** (4.916)	-8.663** (4.374)
sum miss-coordinated outcomes		1.303*** (0.361)
Constant	58.06*** (3.920)	43.78*** (5.216)
Obs.	131	131
$R^2$	0.027	0.072

Standard errors in parentheses  
\*\*\* p<0.01, \*\* p<0.05, \* p<0.1

Table 5: Incentivized guess about the presence of an algorithm in the \_Uncertain treatments, linear probability model.

0.05. In other words, participants are significantly *more* inclined to believe an algorithm was present when this was *not* the case.

**Result 4.** *Guesses about an algorithm being present in the market are significantly lower in Algorithm\_Uncertain\_3 compared to Human\_Uncertain\_3.*

One possible explanation for this surprising finding is that participants associate cooperation with human behavior and not an algorithm.<sup>23</sup> Cooperation rates in Algorithm\_Uncertain are significantly higher than in Human\_Uncertain, and the lower performance of Human\_Uncertain is clearly associated with a higher belief of an algorithm being present. When we add as an explanatory variable to the regressions in Table 5, the number of miscoordinated outcomes (one or two firms chose  $p_{low}$ , whereas at least one firm chose  $p_{high}$ ) which a participant experiences during the entire course of the experiment, this variable is highly significant ( $p < 0.01$ ) and the magnitude of the *algorithm* coefficient decreases, but is still significant ( $p < 0.05$ ). We take this as confirmation that the participants expect the algorithm to be more competitive than humans.

<sup>23</sup>According to Lee (2018), participants rate algorithmic decisions as less fair, trust algorithmic decisions less, and feel less positive about algorithmic decisions when it comes to tasks requiring human skills. With mechanical tasks, the fairness and trustworthiness of algorithms were attributed to their perceived efficiency and objectivity.

### 5.2.3 Differences between human and algorithmic play

One immediate effect of the  $pTFT$  strategy is that it begins a supergame by choosing the high price with probability one, in contrast to the average human subject. Hence, a first attempt at finding differences between humans and the algorithm is to take a closer look at period-one decisions.

Table 6 shows details of regressions similar to those in Section 5.2.1 above, but truncating the data to the first period of each supergame. The algorithm has a substantial and significant effect on the cooperation rate in the first period throughout.<sup>24</sup> This is perhaps not surprising because of the way the algorithm is programmed, but it is important to state this effect formally because of the significance of period-one behavior for overall cooperation.

**Result 5.** *In the  $n = 3$  markets, cooperation rates in the first period are significantly higher in the Algorithm\_ treatments compared to the Human\_ treatments. This also holds when comparing Algorithm\_Uncertain\_3 vs. Human\_Uncertain\_3, and Algorithm\_Certain\_3 vs. Human\_Certain\_3 separately.*

Figure 3 is an alluvial flow diagram that illustrates how humans compare to the algorithm with respect to such individual decisions. It is based on decisions by humans only, using data from all  $n = 3$  treatments, periods 1 to 19 and all supergames.<sup>25</sup> The figure shows how participants' decisions in period  $t - 1$  (left-hand side of the figure) map into market outcomes (middle), and how conditional on these outcomes decisions in period  $t$  emerge. The market outcome is defined as the number of  $p_{high}$  choices of all players in a market, including the subject herself and possibly the algorithm.

---

<sup>24</sup>The effect is also significant comparing the \_Uncertain\_3 and \_Certain\_3 treatments separately. Across all supergames both  $p$ -values  $< 0.01$ .

<sup>25</sup>See in A.7 of the appendix, where we provide the same analysis for the individual treatments. Differences between treatments are minor and insignificant. We dropped the data from period 20 on because we are not specifically interested in the end-game behavior humans exhibit.



	<i>Supergame 1</i>		<i>Supergame 2</i>		<i>Supergame 3</i>		<i>All</i>	
algorithm	0.166*** (0.0446)	0.143** (0.0668)	0.160*** (0.0399)	0.167** (0.0724)	0.202*** (0.0487)	0.225*** (0.0676)	0.176*** (0.0375)	0.178*** (0.0600)
certain	0.0526 (0.0436)	0.0278 (0.0817)	0.0664* (0.0399)	0.0739 (0.0668)	0.0360 (0.0485)	0.0606 (0.0852)	0.0517 (0.0369)	0.0541 (0.0715)
algorithm $\times$ certain		0.0473 (0.0889)		-0.0144 (0.0869)		-0.0468 (0.101)		-0.00462 (0.0789)
supergame							0.0777*** (0.00911)	0.0777*** (0.00896)
Constant	0.321*** (0.0485)	0.333*** (0.0639)	0.444*** (0.0407)	0.440*** (0.0540)	0.465*** (0.0487)	0.453*** (0.0634)	0.332*** (0.0441)	0.331*** (0.0601)
Obs.	309	309	309	309	309	309	927	927
$R^2$	0.031	0.031	0.030	0.030	0.043	0.044	0.050	0.050

Standard errors in parentheses  
\*\*\* p<0.01, \*\* p<0.05, \* p<0.1

Table 6: Cooperation in the first period,  $n = 3$  treatments, linear probability model.

Let us be more specific. Humans choose  $p_{high}$  at a rate of roughly 30% (light blue segment on the left) and, accordingly,  $p_{low}$  at 70% (dark blue segment). Due to the high degree of coordination in markets, outcomes labeled 0 (“all  $p_{low}$ ”) and 3 (“all  $p_{high}$ ”) result most frequently. If the coordination in markets fails, outcomes 1 (“one  $p_{high}$ , two  $p_{low}$ ”) and 2 (“two  $p_{high}$ , one  $p_{low}$ ”) result. The stream from the gray outcome boxes then indicates how humans decided conditional on outcome. Their own  $t - 1$  decision can be identified by the color (light blue for  $p_{high}$  and dark blue for  $p_{low}$ ).

The algorithm always chooses  $p_{high}$  if both competitors previously chose  $p_{high}$ —how do humans behave here? Overall, it turns out human participants are also highly likely to play  $p_{high}$  (92.7%). But there are substantial differences when the own prior choice is taken into account. Provided that they themselves previously played  $p_{high}$ , human subjects almost always play  $p_{high}$  again (99.1%).<sup>26</sup> When we look at the human subjects who played  $p_{low}$  while both their competitors chose  $p_{high}$  (“two  $p_{high}$ , one  $p_{low}$ ”), we see that roughly 29.3% cooperate, whereas the algorithm would play 100%  $p_{high}$  here, too.

Differences between humans and the algorithm also become apparent in markets with mixed outcomes where one competitor chose  $p_{high}$  and the other one  $p_{low}$  in  $t - 1$ . The probability that the algorithm will play cooperatively is 50%, whereas that of the human subjects is only 26.2%. Again, Figure 3 shows the differences between subjects who played  $p_{high}$  previously and those who chose  $p_{low}$ .<sup>27</sup> The cooperatively playing subjects stuck to their strategy with a probability of 61.2%. But such attempts to establish collusive conduct is hampered by the behavior of competitive rivals who rarely choose the high price (9.8%).

How about the potentially negative effect of the algorithm when both rival firms chose  $p_{low}$  previously? In this case, the algorithm would never choose  $p_{high}$ . But this does not differ much for human subjects who cooperate with 3.7%. Conspicuously, the cooperative playing subjects continue

---

<sup>26</sup>In 26 out of 3,007 observations, these subjects chose  $p_{low}$ , which is too little to be visible in Figure 3.

<sup>27</sup>For two-player prisoner’s dilemma experiments, Breitmoser (2015) suggests that subjects play a “semi-grim” strategy, such that subjects randomize across choices regardless of their own previous choice.

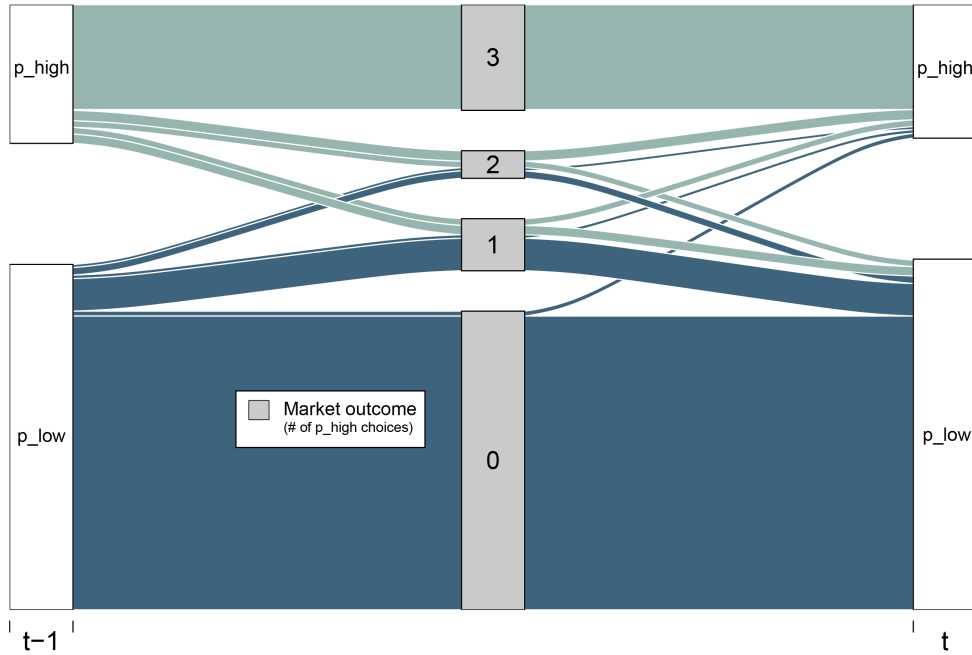


Figure 3: Alluvial flow diagram of choices by human subjects ( $n = 3$  treatments, all supergames, periods 1 to 19).

their strategy with a relatively high probability (41.9%), while the competitive rivals play  $p_{high}$  only in very few cases (1.7%).

Overall, the probability of successful collusion, irrespective of the previous market outcome, is higher in Algorithm\_ (27.4%) than in Human\_ treatments (18.3%). The algorithm is less cooperative than the human subjects when it comes to attempts to establish a collusive outcome, but much more cooperative than subjects who chose  $p_{low}$  before. It seems that the human subjects rarely modify their strategy, trying instead to avoid a change in their price decision.

#### 5.2.4 Profits

If the algorithm variants exhibit more cooperation, this suggests that all firms benefit in terms of higher profits. As we see more cooperation, the mean profits in Algorithm\_ are actually higher than in Human\_, so subjects earn more if an algorithm is present. In Appendix A.3, we analyze this systematically. The positive effect on profits is significant in the third supergame (linear probability model,  $p < 0.05$ ).

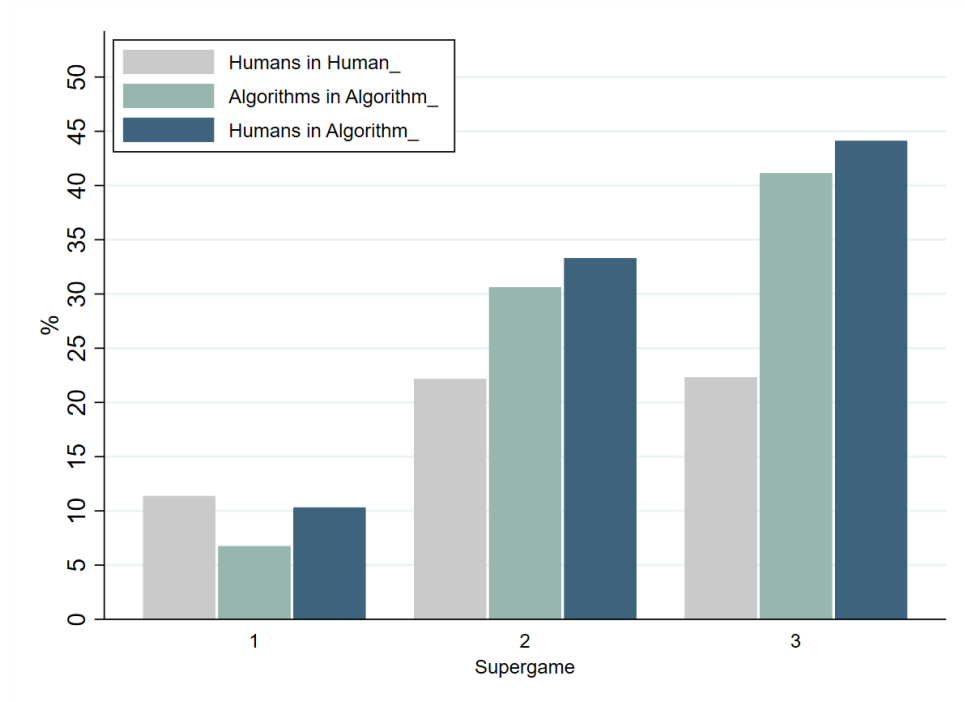


Figure 4: Profits in percent above Nash ( $n = 3$  treatments, periods 6 to 19).

By distinguishing between humans and algorithms, we can analyze who benefits most from the presence of the algorithm. Figure 4 measures profits relative to static Nash earnings (0%) and to perfect collusion (100%).<sup>28</sup> We see that subjects equipped with an algorithm earn substantially less than their competitors in every supergame. Taking all supergames into account, the difference is statistically significant ( $p < 0.01$ ).<sup>29</sup> Although the algorithm helps to increase the group's profit, it performs significantly worse than their competitors. This suggests a coordination problem in that no firm wants to adopt the algorithm first.

**Result 6.** *In the  $n = 3$  markets, profits are significantly higher in the `Algorithm_` treatments compared to the `Human_` treatments. In the `Algorithm_` treatments, participants represented by an algorithm earn significantly less*

<sup>28</sup>Formally, the index in Figure 4 is defined as  $(\pi - \pi^n)/(\pi^c - \pi^n)$ , where  $\pi$  is the observed profit.

<sup>29</sup>Appendix A.3 provides the results of a linear probability regression where the dependent variable is the profit subjects earn from period 6 to 19. Our explanatory variable for the type of player is `role`. The negative effect of `role` is also significant in the first ( $p < 0.1$ ) and the third supergame ( $p < 0.001$ ). The effect is not significant for the second supergame.

than participants who decide themselves for their firm.

### 5.3 Four-firm treatments

We now turn to the experiments with  $n = 4$  firms. The payoffs underlying these sessions can be found in Table 7. Other than that, procedures are virtually identical to the three-firm treatments.

		Other firms' prices			
		<i>All <math>p_{high}</math></i>	<i><math>p_{high}, p_{high}, p_{low}</math></i>	<i><math>p_{high}, p_{low}, p_{low}</math></i>	<i>All <math>p_{low}</math></i>
Own price	<i><math>p_{high}</math></i>	600	0	0	0
	<i><math>p_{low}</math></i>	1,440	720	480	360

Table 7: Payoff table ( $n = 4$  treatments)

Does the collusive effect of the algorithm extend to four-firm markets? On the one hand, the algorithm should promote collusion. On the other hand, we see that already three human subjects are finding it difficult to cooperate. Does it help when we add an algorithm as a fourth player?

Across all supergames, we indeed find more cooperation in Algorithm\_Certain\_4 (0.191) than in Human\_Certain\_4 (0.165). But the difference is small and the effect of the algorithm is not significant, see Appendix A.4. In Figure 2, we see that Algorithm\_Certain\_4 and Human\_Certain\_4 share similarly low weight on outcomes with successful collusion.<sup>30</sup> Consistent with Hypothesis H<sub>5</sub>, we do not see much of a difference between the two treatments with  $n = 4$  firms.

**Result 7.** *Cooperation rates in Algorithm\_Certain\_4 are not significantly higher than in Human\_Certain\_4.*

<sup>30</sup>Subjects equipped with an algorithm still earn less than their competitors. The negative effect of role for the treatments with  $n = 4$  firms is significant in the first ( $p < 0.05$ ) and across all supergames ( $p < 0.1$ ). The effect is not significant for the second and third supergame.

## 6 Conclusion

In this paper, we analyze the impact of algorithms on collusion in hybrid markets where humans interact with algorithms. The analysis of human-computer interaction is important because most markets in the field are heterogeneous and firms cannot be sure of whether their opponents are using algorithms for their pricing decision, nor do they know which type of algorithm competitors might use.<sup>31</sup> A recent and growing literature (Abada and Lambin, 2020; Calvano et al., 2021, 2020b; Klein, 2020) shows that markets with exclusively firms using algorithmic pricing lead to a higher price level. This raises the question of algorithms' impact in hybrid markets where they interact with humans.

We study these issues in experimental markets with three or four firms where one firm is equipped with an algorithm. The algorithm, if present, plays proportional tit-for-tat (Axelrod, 1984; Hilbe et al., 2015). We further vary whether the human participants know (in a non-deceptive way) about the presence of the algorithm. Participants of the experiments played three indefinitely repeated games.

We report three main sets of results. First, regarding the competitiveness of markets, we find that an algorithm significantly increases prices in the three-firm markets. In markets with four firms, the algorithm has no effect on the level of competition. These findings suggest that the collusive impact of algorithms is not a foregone conclusion, but the data likewise indicate the anti-competitive potential algorithms have, even when interacting with humans. Moreover, our findings show that the effects of algorithms facilitating collusion are unlikely to be fully mitigated by the presence of humans. Therefore, we should not rely on humans to discipline the collusive behavior of algorithms.

Our second finding concerns participants' expectations when they interact with an algorithm. Largely, it appears that expectations (the (un)certainty that an algorithm is around) do not significantly affect pricing. Intriguingly, when we elicit post-experiment beliefs about the nature

---

<sup>31</sup>Explicitly communicating and agreeing on the use of algorithms has been penalized as a violation of cartel law. See Poster Cartel case: US Department of Justice, Apr. 6, 2015, Press Release no. 15-421 and British Competition and Markets Authority, Aug. 12, 2016, Case 50223.

of the co-players, participants are significantly more inclined to believe an algorithm was present when this was not the case. Specifically, humans appear to associate miscoordination with algorithmic play whereas, in fact, the algorithm more frequently leads to successful cooperation. These results are broadly consistent with findings on algorithm aversion (Dietvorst and Bharti, 2020; Dietvorst et al., 2016, 2015).

A third set of findings concerns the profitability of employing an algorithm. We find that the firms for which the algorithm decided earn significantly less profit. This finding suggests that firms want their rivals to adopt the algorithm first. In other words, firms face a coordination problem when it comes to delegating decisions to algorithms. Tacit collusion seems feasible, but requires algorithms with a certain degree of cooperative commitment. Therefore, a firm must be willing to accept setup costs. That said, this effect could be moderated by other benefits of algorithms, such as a higher frequency of pricing or a better demand forecasting (Brown and MacKay, 2019; Miklós-Thal and Tucker, 2019).

Our results suggest promising topics for future research. Given that we do not observe higher levels of collusion with four firms, one avenue would be to further explore markets with more firms and a higher number of algorithms. Larger groups may successfully collude if the share or the number of humans is not too high. Another extension would be to not impose the use of the algorithm exogenously, but let subjects choose whether they want to employ algorithms. Algorithm aversion may preclude this, but demonstrating the force of algorithms may cure this reluctance. In addition, the aforementioned coordination problem might be significant.

## References

- Abada, I. and Lambin, X. (2020). Artificial Intelligence: Can Seemingly Collusive Outcomes Be Avoided? *Working Paper*.
- Ahn, T. K., Ostrom, E., Schmidt, D., Shupp, R., and Walker, J. (2001). Cooperation in PD Games: Fear, Greed, and History of Play. *Public Choice*, 106:137–155.
- Assad, S., Clark, R., Ershov, D., and Xu, L. (2020). Algorithmic Pricing and Competition: Empirical Evidence from the German Retail Gasoline Market. *CESifo Working Paper*, No. 8521.
- Axelrod, R. (1984). *The Evolution of Cooperation*. New York: Basic Books.
- Bigoni, M., Casari, M., Skrzypacz, A., and Spagnolo, G. (2015). Time Horizon and Cooperation in Continuous Time. *Econometrica*, 83(2):587–616.
- Blanco, M., Engelmann, D., Koch, A. K., and Normann, H. T. (2014). Preferences and Beliefs in a Sequential Social Dilemma: A Within-Subjects Analysis. *Games and Economic Behavior*, 87:122–135.
- Blonski, M., Ockenfels, P., and Spagnolo, G. (2011). Equilibrium Selection in the Repeated Prisoner’s Dilemma: Axiomatic Approach and Experimental Evidence. *American Economic Journal: Microeconomics*, 3(3):164–192.
- Blonski, M. and Spagnolo, G. (2015). Prisoners’ other Dilemma. *International Journal of Game Theory*, 44:61–81.
- Breitmoser, Y. (2015). Cooperation, But No Reciprocity: Individual Strategies in the Repeated Prisoner’s Dilemma. *American Economic Review*, 105(9):2882–2910.
- British Competition and Markets Authority (2018). Pricing Algorithms: Economic Working Paper on the Use of Algorithms to Facilitate Collusion and Personalised Pricing.



- British Competition and Markets Authority (2021). Algorithms: How They Can Reduce Competition and Harm Consumers.
- Brown, Z. and MacKay, A. (2019). Competition in Pricing Algorithms. *Harvard Business School Working Paper*, no. 20-067.
- Bundeskartellamt and Autorité de la Concurrence (2019). Algorithms and Competition. Discussion Paper.
- Calvano, E., Calzolari, G., Denicolò, V., Harrington, J. E., and Pastorello, S. (2020a). Protecting Consumers from Collusive Prices due to AI: Price-Setting Algorithms can Lead to Noncompetitive Prices, but the Law is Ill Equipped to Stop It. *Science*, 370(6520):1040–1042.
- Calvano, E., Calzolari, G., Denicolò, V., and Pastorello, S. (2020b). Artificial Intelligence, Algorithmic Pricing, and Collusion. *American Economic Review*, 110(10):3267–3297.
- Calvano, E., Calzolari, G., Denicolò, V., and Pastorello, S. (2021). Algorithmic Collusion with Imperfect Monitoring. *International Journal of Industrial Organization (Pre-Proof)*.
- Cameron, A. C., Gelbach, J. B., and Miller, D. L. (2008). Bootstrap-Based Improvements for Inference with Clustered Errors. *Review of Economics and Statistics*, 90(3):414–427.
- Charness, G., Rigotti, L., and Rustichini, A. (2016). Social Surplus Determines Cooperation Rates in the One-Shot Prisoner’s Dilemma. *Games and Economic Behavior*, 100:113–124.
- Chen, L., Mislove, A., and Wilson, C. (2016). An Empirical Analysis of Algorithmic Pricing on Amazon Marketplace. *WWW ’16: Proceedings of the 25th International Conference on World Wide Web*, pages 1339–1349.
- Competition Bureau Canada (2018). Big Data and Innovation: Key Themes for Competition Policy in Canada. Technical report, Competition Bureau Canada.

- Crandall, J. W., Oudah, M., Tennom, Ishowo-Oloko, F., Abdallah, S., Bonnefon, J. F., Cebrian, M., Shariff, A., Goodrich, M. A., and Rahwan, I. (2018). Cooperating with Machines. *Nature Communications*, 9:233.
- Dal Bó, P. and Fréchette, G. R. (2011). The Evolution of Cooperation in Infinitely Repeated Games: Experimental Evidence. *American Economic Review*, 101(1):411–429.
- Dal Bó, P. and Fréchette, G. R. (2018). On the Determinants of Cooperation in Infinitely Repeated Games: A Survey. *Journal of Economic Literature*, 56(1):60–114.
- Dal Bó, P. and Fréchette, G. R. (2019). Strategy Choice in the Infinitely Repeated Prisoner’s Dilemma. *American Economic Review*, 109(11):3929–3952.
- De Melo, C. M., Gratch, J., and Carnevale, P. J. (2015). Humans versus Computers: Impact of Emotion Expressions on People’s Decision Making. *IEEE Transactions on Affective Computing*, 6(2):127–136.
- Degeratu, A. M., Rangaswamy, A., and Wu, J. (2000). Consumer Choice Behavior in Online and Traditional Supermarkets: The Effects of Brand Name, Price, and other Search Attributes. *International Journal of Research in Marketing*, 17(1):55–78.
- Dietvorst, B. J. and Bharti, S. (2020). People Reject Algorithms in Uncertain Decision Domains Because They Have Diminishing Sensitivity to Forecasting Error. *Psychological Science*, 31(10):1302–1314.
- Dietvorst, B. J., Simmons, J. P., and Massey, C. (2015). Algorithm Aversion: People Erroneously Avoid Algorithms after Seeing Them Err. *Journal of Experimental Psychology: General*, 144(1):114–126.
- Dietvorst, B. J., Simmons, J. P., and Massey, C. (2016). Overcoming Algorithm Aversion: People Will Use Imperfect Algorithms If They Can (Even Slightly) Modify Them. *Management Science*, 64(3):1155–1170.
- Dijkstra, J. J., Liebrand, W. B., and Timminga, E. (1998). Persuasiveness of Expert Systems. *Behaviour and Information Technology*, 17:155–163.

- Duffy, J., Hopkins, E., and Kornienko, T. (2021). Facing the Grim Truth: Repeated Prisoner’s Dilemma Against Robot Opponents.
- Duffy, J. and Xie, H. (2016). Group Size and Cooperation Among Strangers. *Journal of Economic Behavior and Organization*, 126:55–74.
- Engel, C. (2015). Tacit Collusion: The Neglected Experimental Evidence. *Journal of Empirical Legal Studies*, 12(3):537–577.
- EU Commission (2017). Final Report on the E-commerce Sector Inquiry. Technical report, European Commission, Brussels.
- Ezrachi, A. and Stucke, M. E. (2016). Virtual Competition. *Journal of European Competition Law & Practice*, 7(9):585–586.
- Ezrachi, A. and Stucke, M. E. (2017). Artificial Intelligence & Collusion: When Computers Inhibit Competition. *University of Illinois Law Review*, 2017(1):1775–1811.
- Farjam, M. and Kirchkamp, O. (2018). Bubbles in Hybrid Markets: How Expectations About Algorithmic Trading Affect Human Trading. *Journal of Economic Behavior and Organization*, 146:248–269.
- Fischbacher, U. (2007). Z-Tree: Zurich Toolbox for Ready-Made Economic Experiments. *Experimental Economics*, 10:171–178.
- Fonseca, M. A. and Normann, H. T. (2012). Explicit vs. Tacit Collusion-The Impact of Communication in Oligopoly Experiments. *European Economic Review*, 56(8):1759–1772.
- Freitag, A., Roux, C., and Thöni, C. (2020). Communication and Market Sharing: An Experiment on the Exchange of Soft and Hard Information. *International Economic Review*, 62(1):175–198.
- Fudenberg, D., Rand, D. G., and Dreber, A. (2012). Slow to Anger and Fast to Forgive: Cooperation in an Uncertain World. *American Economic Review*, 102(2):720–749.
- Green, E. J., Marshall, R. C., and Marx, L. M. (2015). Tacit Collusion in Oligopoly. *The Oxford Handbook of International Antitrust Economics*, 2:464–522.

- Harrington, J. E. (2018). Developing Competition Law for Collusion By Autonomous Artificial Agents. *Journal of Competition Law and Economics*, 14(3):331–363.
- Harrington, J. E. (2020). Third Party Pricing Algorithms and the Intensity of Competition. *Working Paper*.
- Harsanyi, J. C. and Selten, R. (1988). *A General Theory of Equilibrium Selection in Games*.
- Hilbe, C., Wu, B., Traulsen, A., and Nowak, M. A. (2015). Evolutionary Performance of Zero-Determinant Strategies in Multiplayer Games. *Journal of Theoretical Biology*, 374:115–124.
- Honhon, D. and Hyndman, K. (2020). Flexibility and Reputation in Repeated Prisoner’s Dilemma Games. *Management Science*, 66(11):4998–5014.
- Horstmann, N., Krämer, J., and Schnurr, D. (2018). Number Effects and Tacit Collusion in Experimental Oligopolies. *Journal of Industrial Economics*, 66(3):650–700.
- Huck, S., Normann, H. T., and Oechssler, J. (2004). Two are Few and Four Are Many: Number Effects in Experimental Oligopolies. *Journal of Economic Behavior and Organization*, 53(4):435–446.
- Klein, T. (2020). Autonomous Algorithmic Collusion: Q-Learning Under Sequential Pricing. *The RAND Journal of Economics (forthcoming)*.
- Krach, S., Hegel, F., Wrede, B., Sagerer, G., Binkofski, F., and Kircher, T. (2008). Can Machines Think? Interaction and Perspective Taking with Robots Investigated via fMRI. *PLoS ONE*, 3(7):e2597.
- Lee, M. K. (2018). Understanding Perception of Algorithmic Decisions: Fairness, Trust, and Emotion in Response to Algorithmic Management. *Big Data and Society*, 5(1).
- Marwell, G. and Schmitt, D. R. (1972). Cooperation in a Three-Person Prisoner’s Dilemma. *Journal of Personality and Social Psychology*, 21(3):376–383.

- Mehra, S. K. (2016). Antitrust and the Robo-Seller: Competition in the Time of Algorithms. *Minnesota Law Review, Paper No. 2015-15*, 100.
- Miklós-Thal, J. and Tucker, C. (2019). Collusion by Algorithm: Does Better Demand Prediction Facilitate Coordination Between Sellers? *Management Science*, 65(4):1552–1561.
- Monopolkommission (2018). Algorithmen und Kollusion. *Twenty-second Biennial Report by the German Monopolies Commission: Competition 2018*, pages 62–87.
- Murnighan, J. K. and Roth, A. E. (1983). Expecting Continued Play in Prisoner’s Dilemma Games: A Test of Several Models. *Journal of Conflict Resolution*, 27(2):279–300.
- Musolff, L. (2021). Algorithmic Pricing Facilitates Tacit Collusion: Evidence from E-Commerce. *Working Paper*.
- OECD (2017). Algorithms and Collusion: Competition Policy in the Digital Age. Technical report, OECD.
- Osborne, M. J. (2006). Strategic and Extensive Games. *University of Toronto, Department of Economics, Working Papers tecipa-231*.
- Oxera (2017). When Algorithms Set Prices: Winners and Losers. *Oxera Consulting LLP. Discussion Paper*, pages 1–37.
- Potters, J. and Suetens, S. (2013). Oligopoly Experiments in the Current Millennium. *Journal of Economic Surveys*, 27(3):439–460.
- Rilling, J. K., Sanfey, A. G., Aronson, J. A., Nystrom, L. E., and Cohen, J. D. (2004). The Neural Correlates of Theory of Mind Within Interpersonal Interactions. *NeuroImage*, 22(4):1694–1703.
- Romero, J. and Rosokha, Y. (2019). The Evolution of Cooperation: The Role of Costly Strategy Adjustments. *American Economic Journal: Microeconomics*, 11(1):299–328.
- Roth, A. E. and Murnighan, J. K. (1978). Equilibrium Behavior and Repeated Play of the Prisoner’s Dilemma. *Journal of Mathematical Psychology*, 17(2):189–198.

- Roux, C. and Thöni, C. (2015). Collusion Among Many Firms: The Disciplinary Power of Targeted Punishment. *Journal of Economic Behavior and Organization*, 116:83–93.
- Schulz, J., Sunde, U., Thiemann, P., and Thöni, C. (2019). Selection Into Experiments: Evidence from a Population of Students. *Lund University, Department of Economics, Working Papers 2019:18*.
- Waltman, L. and Kaymak, U. (2008). Q-Learning Agents in a Cournot Oligopoly Model. *Journal of Economic Dynamics and Control*, 32(10):3275–3293.
- Weibel, D., Wissmath, B., Habegger, S., Steiner, Y., and Groner, R. (2008). Playing Online Games Against Computer- vs. Human-Controlled Opponents: Effects on Presence, Flow, and Enjoyment. *Computers in Human Behavior*, 24(5):2274–2291.
- Wieting, M. and Sapi, G. (2021). Algorithms in the Marketplace: An Empirical Analysis of Automated Pricing in E-Commerce. *Working Paper*.

# A Appendix

## A.1 Three-player Model

### Setup

Consider a three-player<sup>32</sup> game where players' action sets are the prices  $\{p_{high}, p_{low}\}$ . With  $p_i$  denoting player  $i$ 's price, her payoff is generally denoted by  $\pi_i(p_i, p_j, p_k)$ ,  $i, j, k \in \{1, 2, 3\}$  where  $p_j$  and  $p_k$  are the prices of the rivals of player  $i$ , and  $i \neq j$ ,  $i \neq k$  and  $j \neq k$  and where the identity of the rival players do not matter, that is,  $\pi_i(p_i, p_j, p_k) = \pi_i(p_i, p_k, p_j)$ . We further define the following shortcut notation:

$$\begin{aligned}\pi^c &= \pi_i(p_{high}, p_{high}, p_{high}) = 800 \\ \pi^s &= \pi_i(p_{high}, p_{low}, p_{low}) = \pi_i(p_{high}, p_{low}, p_{high}) = 0 \\ \pi^d &= \pi_i(p_{low}, p_{high}, p_{high}) = 1440 \\ \pi^f &= \pi_i(p_{low}, p_{low}, p_{high}) = 720 \\ \pi^n &= \pi_i(p_{low}, p_{low}, p_{low}) = 480\end{aligned}$$

The numerical entries are those of the experiment and they are also reproduced in Table 1.

We now analyze an infinitely repeated version of this game. Let time be indexed by  $t = 0, \dots, \infty$ . Future periods are discounted by the factor  $\delta$ .

### Repeated-game incentive constraint

Suppose the three players attempt to establish collusion on the high price, each following a 'grim-trigger' strategy ( $GT$ ). When pursuing a  $GT$  strategy, a player chooses  $p_{high}$  in  $t = 0$  and keeps charging  $p_{high}$  as long as no player has played  $p_{low}$  in any previous period. If any player deviates in  $t$ , a  $GT$  player charges  $p_{low}$ , the static Nash equilibrium price, from  $t + 1, \dots, \infty$ . Expected payoffs are as follows. If player  $i$  chooses  $p_{high}$  in  $t = 0$ , she receives  $\pi^c$  from  $t = 0, \dots, \infty$ . If she defects with  $p_{low}$ , she obtains  $\pi^d$  in  $t = 0$  and, since she triggers the punishment path,  $\pi^n$  in  $t = 1, \dots, \infty$ . Accordingly, playing  $GT$  is a subgame-perfect Nash equilibrium (SGPNE)

---

<sup>32</sup>The results below also hold for the  $n = 4$  case. Proof available upon request.

if

$$\begin{aligned} \frac{\pi^c}{1-\delta} &\geq \pi^d + \frac{\delta\pi^n}{1-\delta} \\ \delta &\geq \frac{\pi^d - \pi^c}{\pi^d - \pi^n} = \frac{2}{3} \equiv \underline{\delta}_{GT} \end{aligned} \quad (2)$$

where the subscript  $GT$  indicates that all three participants are  $GT$  players, there is no algorithm.

Suppose now there are two players attempting to establish collusion via  $GT$  and the third player is an algorithm. The algorithm is committed to playing  $pTFT$  (as defined above) and will thus not deviate from this strategy.<sup>33</sup> We analyze the incentives of a  $GT$  player to deviate. If a  $GT$  player chooses  $p_{high}$ , she receives  $\pi^c$  in  $t = 0, \dots, \infty$  in equilibrium. The profit from a one-off deviation is  $\pi^d$ , as before. The punishment payoff in  $t = 1$  does change, however. If player  $i$  deviates in  $t = 0$ , the price vector reads  $(p_{low}, p_{high}, p_{high})$ . This prompts the  $pTFT$  algorithm to cooperate with 50% in  $t = 1$ , so the price vectors  $(p_{low}, p_{low}, p_{high})$  and  $(p_{low}, p_{low}, p_{low})$  and corresponding payoffs  $\pi^f$  and  $\pi^n$ , respectively, are equally likely. From  $t = 2$  on, the two players and the algorithm choose  $p_{low}$  for the rest of the game. Thus, the incentive constraint becomes

$$\frac{\pi^c}{1-\delta} \geq \pi^d + \delta \left( \frac{\pi^f + \pi^n}{2} \right) + \frac{\delta^2\pi^n}{1-\delta}$$

solving for  $\delta$  for the values employed in the experiment<sup>34</sup>

$$\delta \gtrsim 0.69 \equiv \underline{\delta}_{pTFT} \quad (3)$$

where the subscript  $pTFT$  indicates that one of the three players is the  $pTFT$  algorithm.

To complete the proof for the subgame-perfectness of  $GT$  in the presence of the algorithm, consider additional out-of-equilibrium histories.<sup>35</sup> In

---

<sup>33</sup>Tit-for-tat strategies are often not subgame-perfect (for two-player cases, see Osborne (2006)). In our case, the algorithm itself will not deviate, as it is programmed to play  $pTFT$ , even if this is not a best response in some.

<sup>34</sup>A closed-form solution with general payoffs can be obtained, but is not informative.

<sup>35</sup>In the presence of the  $pTFT$  algorithm, meeting the incentive constraint (3) is generally not sufficient for  $GT$  to be subgame-perfect.



histories ending in subgames where at least one player chooses  $p_{low}$  and at least one player selects the high price, the  $GT$  players will choose  $p_{low}$ , whereas the  $pTFT$  algorithm will cooperate in  $t + 1$  with at least 50%.<sup>36</sup> A possible one-off deviation for a  $GT$  player would be to cooperate in the next period. But since the second  $GT$  player will defect in  $t + 1$ , such a deviation would yield zero payoff, whereas sticking to  $GT$  (by defecting from period  $t + 1$  on) would yield (at least)  $\pi^n$ . It follows that  $GT$  is subgame-perfect in the presence of the algorithm, provided (3) is met.

We summarize by comparing (2) and (3):

**Proposition 1:** The minimum discount factor required for collusion to be a SGPNE is lower for three  $GT$  players compared to two  $GT$  players and one  $pTFT$  algorithm:  $\underline{\delta}_{GT} < \underline{\delta}_{pTFT}$ .

The intuition behind Proposition 1 is straightforward.  $GT$  and  $pTFT$  are both cooperative strategies, but  $pTFT$  is more forgiving and willing to cooperate with a positive probability even when (exactly) one rival player defected in  $t - 1$ . This raises the payoffs of a  $GT$  player after a defection and, accordingly, increases the minimum discount factor required for successful collusion.<sup>37</sup>

## Strategic risk

The inequalities (2) and (3) are necessary conditions for collusion on  $p_{high}$  to be subgame-perfect. Other equilibria obviously exist as well. For example, all players always charging  $p_{low}$  is also a SGPNE of the repeated game, with and without the presence of the  $pTFT$  algorithm. The inequalities (2) and (3) do not reflect the coordination problems players face in the presence of multiple equilibria.

---

<sup>36</sup>The set of subgames where at least one player deviates and at least one player cooperates includes  $(p_{high}, p_{low}, \cdot)$ ,  $(p_{low}, p_{high}, \cdot)$  and  $(p_{high}, p_{high}, p_{low})$ . In the latter case, the  $pTFT$  algorithm cooperates 100% in  $t + 1$ . In two further possible subgames ( $(p_{low}, p_{low}, p_{high})$  and  $(p_{low}, p_{low}, p_{low})$ ), all players defect in  $t + 1$ , ensuring  $GT$  is subgame-perfect.

<sup>37</sup>Nevertheless, results from experiments with self-learning algorithms suggest that these algorithms learn to cooperate even after deviations and therefore pursue a more forgiving strategy than  $GT$ , see Calvano et al. (2020b, section IV. C.).

Taking strategic risk into account is especially important when analyzing algorithms. The algorithm is committed to a strategy, whereas humans are not. That is, the algorithm reduces strategic uncertainty. Merely to focus on incentives in a given collusion equilibrium and to ignore strategic risk would imply that we largely miss the collusive impact algorithms may have.

To deal with strategic uncertainty, a growing literature on repeated prisoner’s dilemmas (Blonski et al., 2011; Blonski and Spagnolo, 2015; Dal Bó and Fréchette, 2011, 2018; Green et al., 2015) borrows from Harsanyi and Selten’s (1988) concept of risk dominance which can easily be applied to symmetric coordination games with two strategies. A strategy is risk-dominant if it is a best response to the other players mixing with equal probability between the two strategies.

The approach can be adapted to repeated games. We follow Blonski et al. (2011), Blonski and Spagnolo (2015), Dal Bó and Fréchette (2011, 2018), and Green et al. (2015) in focusing on a simplified version of the game, the choice between two repeated-game strategies. We henceforth analyze the decision between the collusive  $GT$  and the non-cooperative ‘always defect’ strategy ( $AD$ ). That is, players’ action sets are now the repeated-game strategies  $GT$  and  $AD$ .<sup>38</sup> Provided (2) and (3), respectively, hold, all players playing  $GT$  and all players playing  $AD$  are equilibria of this two-action game. Increasing  $\delta$  reduces the riskiness of  $GT$ , and we solve for a new critical discount factor,  $\delta^*$ , such that playing  $GT$  is the best response given the other players randomize with equal probability between the two strategies  $GT$  and  $AD$ . We then investigate how the presence of an algorithm affects  $\delta^*$ .

Consider three players choosing between  $GT$  and  $AD$  and expecting their competitors to play  $GT$  or  $AD$  with equal probability. When playing  $GT$ , there are two contingencies for the profit of player  $i$  in period  $t = 0$ : Provided the other two players also play  $GT$  (which happens with a probability of  $1/4$ ),  $i$  obtains  $\pi^c$ . If at least one other player defects (probability of  $3/4$ ),  $i$  obtains  $\pi^s = 0$  in period  $t = 0$ . If all players including  $i$  play  $GT$

---

<sup>38</sup>For the simplified version of the game with only two repeated-game strategies ( $GT$  and  $AD$ ), Blonski and Spagnolo (2015) show that any collusive equilibrium is risk-dominant if  $GT$  is risk-dominant.

in  $t = 0$ ,  $i$  also obtains  $\pi^c$  in all future periods  $t = 1, \dots, \infty$ . If at least one player defects in  $t = 0$ ,  $i$  gets  $\pi^n$  in periods  $t = 1, \dots, \infty$ . Thus, player  $i$ 's expected payoff from playing  $GT$  is

$$\frac{1}{4} \left( \frac{\pi^c}{1 - \delta} \right) + \frac{3}{4} \left( \pi^s + \frac{\delta \pi^n}{1 - \delta} \right)$$

If player  $i$  instead plays  $AD$ , there are three possibilities. If both other players cooperate in  $t = 0$  (which happens with a probability of  $1/4$ ),  $i$  obtains  $\pi^d$ . If one rival player cooperates and the other defects (which happens with a probability of  $1/2$ ),  $i$  obtains  $\pi^f$ . When both rival players defect (probability of  $1/4$ ),  $i$  obtains  $\pi^n$ . In all three cases,  $i$  obtains  $\pi^n$  in  $t = 1, \dots, \infty$ . Player  $i$ 's expected payoff is

$$\frac{\pi^d}{4} + \frac{\pi^f}{2} + \frac{\pi^n}{4} + \frac{\delta \pi^n}{1 - \delta}$$

Taking the difference in expected profits of  $GT$  and  $AD$  and solving for  $\delta$ , we find that  $GT$  has a higher expected payoff than  $AD$ , if and only if

$$\delta \geq \frac{\pi^d + 2\pi^f - 3\pi^s + \pi^n - \pi^c}{\pi^d + 2\pi^f - 3\pi^s} = \frac{8}{9} \approx 0.89 \equiv \delta_{GT}^* \quad (4)$$

where  $\delta^* \in (0, 1)$  denotes the critical discount factor in the presence of strategic risk and  $\delta^*$ .  $GT$  indicates that all three players are (potential)  $GT$  players. Note that  $\delta_{GT}^* > \underline{\delta}_{GT}$  strictly. We further point out that the payoffs  $\pi^s$  and  $\pi^f$  play a role here – which is not the case for  $\underline{\delta}$ .

Now, one of the three market participants is an algorithm committed to playing  $pTFT$ , whereas the other two participants are rational players as before. We analyze the choice of these two players between  $GT$  and  $AD$ . The two players are expected to play  $GT$  and  $AD$  with equal probability. Suppose player  $i$  plays  $GT$ . Then there are only two contingencies: the other player plays either  $GT$  or she plays  $AD$ . Expected profits are accordingly

$$\frac{1}{2} \left( \frac{\pi^c}{1 - \delta} \right) + \frac{1}{2} \left( \pi^s + \frac{\delta}{2} (\pi^f + \pi^n) + \delta^2 \frac{\pi^n}{1 - \delta} \right)$$

If player  $i$  plays  $AD$ , she gets

$$\frac{1}{2} \left( \pi^d + \frac{\delta}{2} (\pi^f + \pi^n) + \frac{\delta^2 \pi^n}{1 - \delta} \right) + \frac{1}{2} \left( \pi^f + \frac{\delta \pi^n}{1 - \delta} \right)$$

We find that  $GT$  has a higher expected payoff if

$$\delta \geq \frac{\pi^d + \pi^f - \pi^c - \pi^s}{\pi^d + \pi^f - \pi^n - \pi^s} = \frac{17}{21} \approx 0.81 \equiv \delta_{pTFT}^* \quad (5)$$

Comparing (4) and (5), we obtain:

**Proposition 2:** The minimum discount factor required in the presence of strategic uncertainty is higher for three  $GT$  players compared to two  $GT$  players and one  $pTFT$  algorithm:  $\delta_{GT}^* > \delta_{pTFT}^*$ .<sup>39</sup>

Whereas propositions 1 and 2 imply contradicting effects, the existing experimental evidence overwhelmingly suggests that the minimum discount factor that takes strategic risk into account ( $\delta^*$ ) has more explanatory power than the standard minimum discount factor ( $\underline{\delta}$ ). This is shown in the meta-study by Dal Bó and Fréchet (2018). Furthermore, Blonski et al. (2011) highlight the case where, for pairs of treatments,  $\underline{\delta}$  and  $\delta^*$  “change in opposite directions.” This is the case in our experiment: the  $pTFT$  algorithm increases  $\underline{\delta}$ , but reduces  $\delta^*$ . The first experimental result in Blonski et al. (2011) is that, in this case, “the frequency of cooperation changes as predicted by changes in  $\delta^*$ , contradicting predictions based on  $\underline{\delta}$ ” (Blonski et al., 2011, p. 185).<sup>40</sup> Accordingly, based on Proposition 2, we expect that, ceteris paribus, the Algorithm\_ treatments will be more collusive than their Human\_ counterparts.

---

<sup>39</sup>The reader can verify that  $\delta_{GT}^* > \delta_{pTFT}^*$  not only for our experimental parameters, but in general: Note that both the numerator and the denominator of  $\delta_{GT}^*$  exceed their  $\delta_{pTFT}^*$  counterparts by  $\pi^f + \pi^n - 2\pi^s > 0$ , hence are increasing  $\delta_{GT}^*$ .

<sup>40</sup>Their analysis is often based on what they label as “class 2” data. In that class, the actual discount factor is above  $\underline{\delta}$ , but below  $\delta^*$ , as is the case in our experiment.

## A.2 Treatment ranking as explanatory variable

	<i>Supergame 1</i>	<i>Supergame 2</i>	<i>Supergame 3</i>	<i>All</i>
treatment ranking	0.00432 (0.0217)	-0.0556*** (0.0203)	-0.0879** (0.0342)	-0.0439** (0.0191)
periods 1 to 5	0.0867*** (0.0234)	0.129*** (0.0177)	0.116*** (0.0199)	0.110*** (0.0133)
periods 20 to 25	-0.0524*** (0.0159)	-0.0522*** (0.0160)	-0.129*** (0.0292)	-0.0821*** (0.0108)
supergame				0.0968*** (0.0216)
Constant	0.145** (0.0568)	0.438*** (0.0447)	0.566*** (0.103)	0.280*** (0.0521)
Obs.	7,416	6,180	6,489	20,085
$R^2$	0.016	0.033	0.062	0.065

Standard errors in parentheses

\*\*\* p<0.01, \*\* p<0.05, \* p<0.1

Table 8: Impact of treatment ranking,  $n = 3$ , periods 6 to 19, linear probability model.

### A.3 Profits

	<i>Supergame 1</i>		<i>Supergame 2</i>		<i>Supergame 3</i>		<i>All</i>	
algorithm	62.63 (274.8)	116.0 (274.6)	832.5*** (300.2)	872.5*** (303.6)	1,348** (600.2)	1,393** (599.1)	714.8*** (265.9)	761.2*** (266.5)
certain	491.2 (539.1)	491.2 (539.1)	249.6 (385.1)	249.6 (385.1)	289.1 (591.3)	289.1 (591.3)	351.6 (402.5)	351.6 (402.5)
algorithm $\times$ certain	-330.3 (633.9)	-330.3 (633.9)	-770.9* (454.0)	-770.9* (454.0)	-857.9 (905.8)	-857.9 (905.8)	-636.3 (538.5)	-636.3 (538.5)
role		-160* (91.29)		-120 (78.87)		-133.3*** (14.63)		-139.1*** (50.78)
supergame							521.1*** (106.7)	521.1*** (106.7)
Constant	6,989*** (188.7)	6,989*** (188.7)	7,590*** (252.8)	7,590*** (252.8)	7,578*** (285.9)	7,578*** (285.9)	6,867*** (112.3)	6,867*** (112.3)
Obs.	7,416	7,416	6,180	6,180	6,489	6,489	20,085	20,085
$R^2$	0.014	0.016	0.022	0.023	0.059	0.059	0.065	0.066

Standard errors in parentheses

\*\*\* p<0.01, \*\* p<0.05, \* p<0.1

Table 9: Total profits,  $n = 3$  treatments, periods 6 to 19, linear probability model.

#### A.4 Treatment effects for $n = 4$ variants

	<i>Supergame 1</i>	<i>Supergame 2</i>	<i>Supergame 3</i>	<i>All</i>
algorithm	0.0725 (0.0873)	0.0117 (0.118)	0.00604 (0.0951)	0.0296 (0.0910)
periods 1 to 5	0.136*** (0.0206)	0.160*** (0.0368)	0.164*** (0.0316)	0.153*** (0.0170)
periods 20 to 25	-0.0455* (0.0266)	-0.153*** (0.0526)	-0.0855*** (0.0266)	-0.101*** (0.0247)
supergame				0.0413*** (0.0158)
Constant	0.0839*** (0.0201)	0.206** (0.0891)	0.200*** (0.0490)	0.122*** (0.0418)
Obs.	2,904	3,300	2,772	8,976
$R^2$	0.042	0.066	0.034	0.051

Standard errors in parentheses

\*\*\*  $p < 0.01$ , \*\*  $p < 0.05$ , \*  $p < 0.1$

Table 10: Treatment effects,  $n = 4$  variants, all periods, linear probability model.

# Additional Material

## A.5 Session details

Date	Session	# Part.	# Markets	Lab <sup>a</sup>	COVID-19
August 28, 2019	31	21	7	Düsseldorf	0
September 11, 2019	41	24	8	Düsseldorf	0
September 11, 2019	32	21	7	Düsseldorf	0
September 12, 2019	42	18	6	Düsseldorf	0
February 19, 2020	51	16	4	Düsseldorf	0
February 19, 2020	52	20	5	Düsseldorf	0
March 04, 2020	33	24	8	Bonn	0
March 05, 2020	11	21	7	Bonn	0
March 05, 2020	12	21	7	Bonn	0
March 05, 2020	21	30	10	Bonn	0
July 06, 2020	22	18	6	Düsseldorf	1
July 15, 2020	61	20	5	Düsseldorf	1
August 05, 2020	13	18	6	Düsseldorf	1
August 17, 2020	62	20	5	Düsseldorf	1
September 02, 2020	23	15	5	Düsseldorf	1
September 22, 2020	43	18	6	Bonn	1
September 22, 2020	63	12	3	Bonn	1
September 22, 2020	53	16	4	Bonn	1
October 13, 2020	44	18	6	Bonn	1
October 13, 2020	64	12	3	Bonn	1
October 13, 2020	54	16	4	Bonn	1
October 14, 2020	24	12	4	Bonn	1
October 14, 2020	34	18	6	Bonn	1
October 16, 2020	14	12	4	Düsseldorf	1

Table 11: Session details

<sup>a</sup>“Corona” indicates that the session was conducted under the common hygiene rules of the pandemic. As a show-up fee, the participants received 5 euros in Bonn and 4 euros in Düsseldorf. During the COVID-19 pandemic, the fee was increased to 8 euros in Düsseldorf from mid-July 2020 on. This in line with Schulz et al. (2019), who find that moderately different show-up fees had no influence on the behavior of the participants.



## A.6 Overview using data from all periods

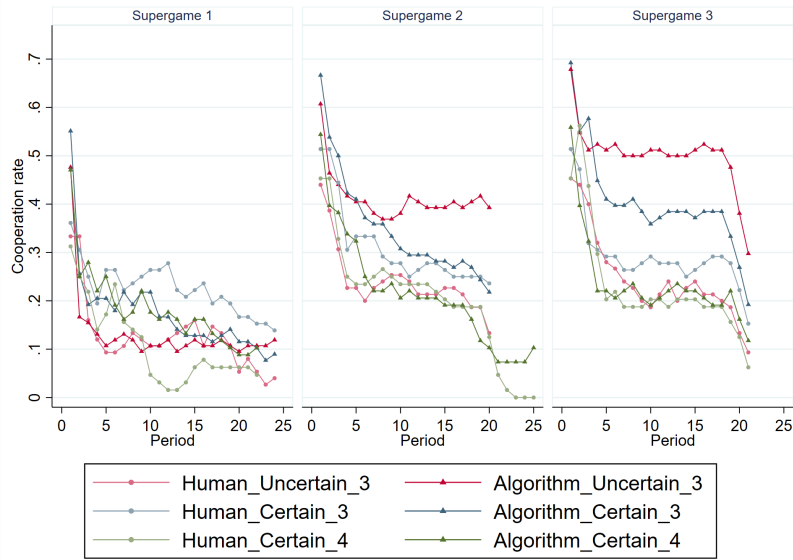


Figure 5: Cooperation rates over time (all periods).

	<i>Supergame 1</i>	<i>Supergame 2</i>	<i>Supergame 3</i>	<i>All</i>
Human_U_3	0.126 (0.331)	0.241 (0.428)	0.246 (0.431)	0.200 (0.400)
Algo_U_3	0.130 (0.337)	0.412 (0.492)	0.502 (0.500)	0.337 (0.473)
Human_C_3	0.226 (0.418)	0.310 (0.463)	0.293 (0.455)	0.273 (0.446)
Algo_C_3	0.174 (0.379)	0.350 (0.477)	0.404 (0.491)	0.302 (0.459)
Human_C_4	0.109 (0.311)	0.201 (0.401)	0.231 (0.421)	0.180 (0.384)
Algo_C_4	0.181 (0.385)	0.212 (0.409)	0.237 (0.425)	0.210 (0.407)

Standard deviations in parentheses.

Table 12: Average cooperation rates (all periods).

## A.7 Data of human subjects only

	<i>Supergame 1</i>	<i>Supergame 2</i>	<i>Supergame 3</i>	<i>All</i>
Human_U_3	0.123 (0.328)	0.221 (0.415)	0.218 (0.413)	0.187 (0.390)
Algorithm_U_3	0.112 (0.316)	0.395 (0.489)	0.504 (0.500)	0.337 (0.473)
Human_C_3	0.233 (0.423)	0.275 (0.447)	0.277 (0.448)	0.262 (0.440)
Algorithm_C_3	0.159 (0.366)	0.297 (0.457)	0.378 (0.485)	0.278 (0.448)
Human_C_4	0.0804 (0.272)	0.222 (0.416)	0.193 (0.395)	0.165 (0.371)
Algorithm_C_4	0.155 (0.363)	0.197 (0.398)	0.214 (0.411)	0.189 (0.392)

Standard deviations in parentheses.

Table 13: Average cooperation rates (human subjects, periods 6 to 19).

	Rival behavior in $t - 1$			
	<i>Two Low</i>	<i>High/Low</i>	<i>Two High</i>	<i>Total</i>
Human_U_3	0.0401 (0.196)	0.267 (0.443)	0.908 (0.289)	0.207 (0.405)
Algorithm_U_3	0.0291 (0.168)	0.238 (0.427)	0.953 (0.211)	0.343 (0.475)
Human_C_3	0.0405 (0.197)	0.258 (0.438)	0.948 (0.221)	0.278 (0.448)
Algorithm_C_3	0.0347 (0.183)	0.281 (0.450)	0.886 (0.318)	0.301 (0.459)
Total	0.0370 (0.189)	0.262 (0.440)	0.927 (0.260)	0.276 (0.447)

Standard deviations in parentheses.

Table 14: Average cooperation rates with respect to the previous choices of rivals 1 and 2 (human subjects,  $n = 3$  treatments, period 6-19).

## A.8 Treatments effect with probit model

	<i>Supergame 1</i>	<i>Supergame 2</i>	<i>Supergame 3</i>	<i>Supergame 3</i>	<i>Supergame 3</i>	<i>Supergame 3</i>	<i>All</i>	<i>All</i>
algorithm	-0.0907 (0.223)	0.0284 (0.293)	0.304* (0.156)	0.490** (0.203)	0.509** (0.235)	0.704** (0.311)	0.258 (0.169)	0.436** (0.172)
certain	0.294 (0.224)	0.406 (0.356)	0.00233 (0.149)	0.210 (0.258)	-0.0760 (0.237)	0.147 (0.360)	0.0610 (0.166)	0.257 (0.274)
algorithm × certain		-0.220 (0.491)		-0.378 (0.296)		-0.400 (0.501)		-0.357 (0.337)
periods 1 to 5	0.320*** (0.0968)	0.321*** (0.0973)	0.349*** (0.0493)	0.350*** (0.0492)	0.309*** (0.0578)	0.311*** (0.0572)	0.320*** (0.0420)	0.322*** (0.0410)
periods 20 to 25	-0.252*** (0.0841)	-0.254*** (0.0841)	-0.160*** (0.0555)	-0.163*** (0.0562)	-0.395*** (0.0798)	-0.398*** (0.0806)	-0.322*** (0.0483)	-0.325*** (0.0476)
supergame							0.300*** (0.0647)	0.301*** (0.0633)
Constant	-1.122*** (0.175)	-1.184*** (0.169)	-0.687*** (0.159)	-0.792*** (0.185)	-0.625*** (0.186)	-0.737*** (0.183)	-1.103*** (0.147)	-1.206*** (0.119)
Obs.	7,416	7,416	6,180	6,180	6,489	6,489	20,085	20,085

Standard errors in parentheses

\*\*\* p<0.01, \*\* p<0.05, \* p<0.1

Table 15: Treatment effects,  $n = 3$  variants, probit model.

## A.9 Robustness check for impact of lab location and COVID-19 pandemic

	<i>Supergame 1</i>		<i>Supergame 2</i>		<i>Supergame 3</i>		<i>All</i>	
algorithm	-0.0225 (0.0528)	-0.00359 (0.107)	0.108** (0.0526)	0.178* (0.106)	0.185** (0.0813)	0.255 (0.196)	0.0847 (0.0516)	0.136 (0.109)
certain	0.0706 (0.0542)	0.0346 (0.114)	0.000116 (0.0515)	0.0712 (0.138)	-0.0292 (0.0833)	0.0457 (0.219)	0.0166 (0.0523)	0.0494 (0.132)
algorithm × certain		0.0174 (0.169)		-0.138 (0.168)		-0.142 (0.294)		-0.0821 (0.176)
corona <sup>a</sup>		-0.0277 (0.0928)		0.0177 (0.0823)		-0.00476 (0.165)		-0.00634 (0.0881)
laboratory <sup>b</sup>		0.0368 (0.0747)		-0.0110 (0.0733)		0.0182 (0.131)		0.0161 (0.0730)
periods 1 to 5	0.0867*** (0.0234)	0.0896*** (0.0236)	0.129*** (0.0177)	0.124*** (0.0175)	0.116*** (0.0199)	0.116*** (0.0210)	0.110*** (0.0133)	0.110*** (0.0135)
periods 20 to 25	-0.0524*** (0.0159)	-0.0512*** (0.0161)	-0.0522*** (0.0160)	-0.0534*** (0.0176)	-0.129*** (0.0292)	-0.134*** (0.0306)	-0.0821*** (0.0108)	-0.0842*** (0.0114)
supergame							0.0968*** (0.0216)	0.105*** (0.0206)
Constant	0.133*** (0.0405)	0.114 (0.106)	0.245*** (0.0518)	0.208** (0.102)	0.268*** (0.0625)	0.224 (0.165)	0.120*** (0.0401)	0.0795 (0.101)
Obs.	7,416	7,128	6,180	5,940	6,489	6,237	20,085	19,305
$R^2$	0.025	0.025	0.029	0.034	0.057	0.064	0.062	0.074

Standard errors in parentheses  
\*\*\* p<0.01, \*\* p<0.05, \* p<0.1

Table 16: Laboratory location and COVID-19 effects, n = 3 variants, linear probability model.

<sup>a</sup>corona = 1 if the session was conducted under hygiene rules of the pandemic.

<sup>b</sup>laboratory = 1 if the sessions was run in Bonn.

## **A.10 Instructions (\_Uncertain treatments)**

### **Welcome to the experiment**

Thank you for your participation in this experiment. Please read the instructions carefully. For your participation in today's experiment, you will receive 5 euros. During the experiment, you will have the opportunity to earn an additional amount of money. The additional amount will depend on your decisions and the decisions of the other participants. A short questionnaire will follow the experiment. From now on, please stop any conversations with your neighbors. Turn off your cell phone and remove everything from your table that you do not need for the experiment. If you have any questions, please raise your hand and we will answer them one-on-one.

### **Instructions**

In this experiment, you will take the role of a firm in a market. Each market consists of three firms. Each of the three firms is represented by a human participant. All firms offer 24 units of a comparable good with no cost of production, and with 24 consumers demanding one unit of the good. Consumers' willingness to pay for a good ranges from 1 to 100 ECU (Experimental Currency Units), where 1,000 ECU = 1 Euro. At the beginning of each period, all firms have the option to set a high price (100 ECU) or a low price (60 ECU) for their good. The company which alone has set the lowest price serves the entire demand. All other companies will not sell any of their units. If several companies have set the same lowest price, the demand is divided equally among them. The following three examples illustrate the mechanism of the market:

#### **Example 1**

You are firm A and you decide to charge a high price for the units of your good (100 ECU). Firm B makes the same decision, whereas C sets a low price (60 ECU). Firm C now has the cheapest sales offer and will serve the complete demand. Accordingly, firm C will earn (60 ECU \* 24 units sold =) 1,440 ECU. Firms A and B will not sell any units and will therefore

	Both competitors choose the high price	One competitor chooses the high price, the other competitor chooses the low price	Both competitors choose the low price
<b>You choose the high price (100 ECU)</b>	800 ECU	0 ECU	0 ECU
<b>You choose the low price (60 ECU)</b>	1440 ECU	720 ECU	480 ECU

earn 0 ECU in this period.

### **Example 2**

You are firm A and you decide to charge a low price for the units of your good (60 ECU). Firms B and C make the same decision. Firms A, B, and C have now all made the lowest sales offer and will each sell  $1/3$  of the demand, thus  $24/3 = 8$  units of their goods. Accordingly, each firm will earn (60 ECU \* 8 units sold =) 480 ECU.

### **Example 3**

You are firm A and you decide to charge a high price for the units of your good (100 ECU). Firms B and C make the same decision. Firms A, B, and C have now all made the most favorable sales offer and will each sell  $1/3$  of the demand, thus  $24/3 = 8$  units of their goods. Accordingly, each firm will earn (100 ECU \* 8 units sold =) 800 ECU. Thus, your earnings depend on your own and the other firms' pricing decisions. This results in the following profit table for you:

After all the firms have made their choice, you will be informed about the chosen prices of the other two firms and about your profit.

### **Periods and rounds**

In total, you will play at least 20 periods with the other two firms. Random chance will decide whether or not additional periods will be played in the sequel. With a probability of 70% the round will continue with another

period; with a probability of 30% the round will end. The round continues until random chance determines the end. In each period of a round, you will be playing with the same participants in a market. At the end of these  $20 + x$  periods, all participants are randomly assigned to new markets and a new round begins. The three participants in the new markets will then stay together again for  $20 + x$  periods.

In total, you will play three rounds of  $20+x$  periods. After three rounds, the experiment ends and a short questionnaire follows.

### **Market decisions by algorithms**

In your markets, at least two participants decide for themselves the price for which they want to sell their goods for their firm and are paid the profit their firm makes in cash at the end of the experiment. **With 50% probability, the decisions for the third firm will also be made by one participant. Also with 50% probability, the third firm will be equipped with an algorithm in all rounds, which will make the necessary pricing decisions for the participant.** In this case, the participant does not make any decisions but still gets paid in cash the profit that her firm makes.

### **Payout**

For your payout, one of the three rounds will be randomly selected. The ECU earned there will be paid to you additionally in euros.