

WP 2302 – January 2023

An Experimental Test of Algorithmic Dismissals

Brice Corgnet

Abstract:

We design a laboratory experiment in which a human or an algorithm decides which of two workers to dismiss. The algorithm automatically dismisses the least productive worker whereas human bosses have full discretion over their decisions. Using performance metrics and questionnaires, we find that fired workers react more negatively to human than to algorithmic decisions in a broad range of tasks. We show that spitefulness exacerbated this negative reaction. Our findings suggest algorithms could help tame negative reactions to dismissals.

Keywords:

Algorithmic dismissals, laboratory experiments, distributive justice, work satisfaction, social preferences

JEL codes:

C92, D23, D91, M50, O33

An Experimental Test of Algorithmic Dismissals

Brice Corgnet[†]

Abstract

We design a laboratory experiment in which a human or an algorithm decides which of two workers to dismiss. The algorithm automatically dismisses the least productive worker whereas human bosses have full discretion over their decisions. Using performance metrics and questionnaires, we find that fired workers react more negatively to human than to algorithmic decisions in a broad range of tasks. We show that spitefulness exacerbated this negative reaction. Our findings suggest algorithms could help tame negative reactions to dismissals.

Keywords: Algorithmic dismissals, laboratory experiments, distributive justice, work satisfaction, social preferences.

JEL codes: C92, D23, D91, M50, O33

1. Introduction

Humans have long feared being replaced by machines (e.g., Keynes, 1930). New technological progress in artificial intelligence has intensified this threat to the point that Stuart Russell (2019) claimed that: “Success in creating AI would be the biggest event in human history. Unfortunately, it might also be the last.”¹ The advent of AI could render the expertise of numerous professionals obsolete (e.g., Brynjolfsson and McAfee, 2014; Frey and Osborne, 2017; Daugherty and Wilson, 2018; Manyika and Sneider, 2018), thus confirming the conjecture of Simon (1965). A novel feature of AI-based technological progress is the rising level of autonomy of machines that can teach themselves as in some applications of deep learning techniques (Russell and Norvig, 2021).

[†] Emlyon Business School, GATE UMR 5824, F-69130 Ecully, France. Email: corgnet@emlyon.com. Brice acknowledges Quentin Thévenet for his help in programming and conducting the experiment. Many thanks to the Economic Science Institute at Chapman University for their lifelong support.

¹ This quote is often attributed to Stephen Hawking et al., (2014) but Russell (2019, p.4) points to his 2013 Dulwich lecture.

It follows that AI-powered machines will be able to make critical decisions without the intervention of humans. This is deeply troubling because decision-making, especially when it is deliberate and reflective (Stanovich, 2005), is often considered a key dimension of what defines us as humans (Gazzaniga, 2000).

Deliberate thinking is part and parcel of many jobs, especially those requiring management skills. Among the most challenging decisions managers face, are those related to employees' evaluations and dismissals (Badaracco, 1997, 2016). Yet, companies are starting to use algorithms to assist or even replace human managers. In a recent case, *Amazon* has admitted to dismissing warehouse workers based on algorithmic evaluations without recurring to any human judgment (Lecher, 2019). They motivated the use of algorithms by their willingness to avoid human biases in dismissal decisions, thus focusing exclusively on productivity metrics. Hoffman, Kahn and Li (2018) show that managers who overrule test metrics when hiring tend to recruit worse candidates, and interpret this finding as resulting from human biases. Algorithmic dismissals are also particularly cheap, and this should be no surprise that other tech companies, such as Uber or Xsolla (see Smith, 2020; Fortson, 2021), are deploying similar tools.

However, the *Amazon* legal case had unintended consequences because algorithmic dismissals were largely condemned in the media. These concerns are exemplified by the recent rules established by European regulators to restrain the use of algorithmic decisions (GDPR, General Data Protection Regulation). Of particular interest is Article 22(1) according to which "A person shall have the right not to be subject to a decision based solely on automated processing (...)." Scholars and practitioners have stressed the ambiguity of this statement, thus casting doubt on its effective implementation (see Russell, 2019). Daugherty and Wilson (2018, p.124) and Russell (2019, p.128) interpret GDPR rules as a "right to an explanation" for automated decisions, thus leaving open the possibility to use algorithmic dismissals.

Not only the legal system lags technological advances but also the research on algorithmic decision-making. We are not aware of any scientific study of algorithmic dismissals. In this paper, we start filling this gap by studying the reaction of the workers who are directly impacted by automated dismissals. Doing so, we leave aside other relevant dimensions such as the evaluation of these practices by the general public as well as the legal issues associated with algorithmic dismissals. Specifically, we ask two related questions: 1) Will algorithmic dismissals hamper or

foster the performance of human workers compared to human dismissals? 2) Will they hamper or foster workers' perception of distributive justice and satisfaction on the job? Answering these two questions will help us assess whether algorithmic dismissals can be effectively implemented in an organization. Indeed, effective implementation will not be possible if automated dismissals lead to lower productivity, exacerbated feelings of injustice and job dissatisfaction.

The lack of research on algorithmic dismissals can be explained by the unavailability of archival data. Access to this data is limited because these techniques are controversial, and in some cases even illegal given GDPR regulations. Furthermore, companies are still experimenting with these techniques and are thus reluctant to share any data that could threaten their competitive advantage. A promising approach is to collect experimental data.² In this paper, we propose a new experimental paradigm to compare human and algorithmic dismissals. In our experiment, a boss and two workers interact across 7 periods that are divided in two stages. In the first stage of each period, the two workers receive a 5-euro fixed pay to complete a task that consists in reproducing patterns using pixels. Before the second stage starts, the boss decides which employee to maintain at the workstation after observing feedback about the production of both workers. The maintained workers keep their first-stage pay of 5 euros whereas demoted workers are only rewarded 1 euro for completing the task.

We use a between-subject design with two treatments. In the *Human* treatment, the boss is a human participant whereas in the *Algo* treatment an algorithm makes the decision. The algorithm only uses performance metrics and dismisses the least productive worker in the first period. This simple rule aims at replicating the basic features of dismissal algorithms currently used by companies like *Amazon* that focus on performance alone. Because this rule is easy to describe to workers, it allows us to alleviate common transparency concerns associated with algorithmic decisions (Shin and Park, 2019; De Cremer, 2020). By the same token, our simple algorithm minimizes issues related to perceived biases in automated decisions, which is the subject of a growing literature (see Cowgill and Tucker, 2020). These design choices allow us to isolate *intentionality* as a distinctive feature of human versus algorithmic decisions (Hidalgo, 2021). Intentionality is a key dimension of human decisions that, unlike transparency and biases, cannot be replicated by the machine.

² See Chuginova and Sele (2020) for a review of experimental papers on human-machine interactions.

In our design, we also vary the nature of the task across periods to identify job characteristics that impact workers' reactions to algorithmic dismissals. In the case of *measurable* tasks, algorithms and human bosses can calculate performance without error. In that case, the only difference between human and algorithmic dismissal decisions relates to intentionality. In the case in which task performance is difficult to measure (*non-measurable* tasks), algorithms have an edge over human bosses because they can compute productivity metrics without error. We also consider tasks in which human bosses would tend to favor one worker based on a dimension unrelated to task performance (*bias* tasks). Finally, we consider cases in which one of the two workers had a more difficult task to complete, thus facing a handicap (*handicap* tasks). In that case, human bosses could decide, unlike algorithms, to disregard performance metrics and reward the effort of the handicapped worker.

Building on the literature in behavioral economics and human-robot interaction, we posit four hypotheses. In the first three hypotheses, we conjecture that fired workers will react more positively to algorithmic than to human decisions in a broad range of tasks. This is the case for easy and measurable tasks in which workers obtain a similar level of performance but one of them must be fired (Hypothesis *Ii*). In that case, human dismissals exacerbate the perception of injustice because they are seen as intentional. We expect the difference between human and algorithmic decisions to be less pronounced for hard and measurable tasks (Hypothesis *Iii*). In that case, workers will observe stark differences in performance between themselves so that human dismissals that are based on relative performance will likely be perceived as fair. In the case of bias and handicap tasks (Hypotheses 2 and 3), dismissals will be perceived as unfair, thus triggering more negative reactions in the *Human* than in the *Algo* treatment. Finally, in Hypothesis 4, we conjecture that the difference in attitudes towards human and algorithmic dismissals can be explained by workers' levels of algorithm aversion (Syrdal et al., 2009).

Our experimental results provide evidence for the first three hypotheses. Fired workers reacted more positively than non-fired workers when algorithmic rather than human dismissals were used. Overall, workers reacted more negatively when fired by a human than when fired by an algorithm. However, this pattern depended on the nature of the task. For a measurable task in which one of the two workers clearly outperformed the other, fired workers reacted similarly to human and algorithmic dismissals in line with Hypothesis *Iii*.

Overall, the effect of algorithmic dismissals on workers' performance was positive. This result might appear surprising given that we are overwhelmed with anecdotes and opinion pieces emphasizing the negative reaction of workers to algorithmic dismissals (Lecher, 2019).

On a cautionary note, non-fired workers perceived being picked by a human to be fairer than being selected by an algorithm. As a result, they tended to perform better after being picked by a human boss than after being chosen by an algorithm, but this difference was not significant for most tasks. Finally, in contrast with Hypothesis 4, we found no evidence that workers' level of algorithm aversion could explain their reaction to algorithmic dismissals.

2. Related literature and hypotheses

2.1. Algorithmic decisions

2.1.1. Supervision and evaluation

Our study of algorithmic dismissals naturally connects to the literature on automated supervision and work evaluations (e.g., Raveendhran and Fast, 2021; Fumagalli, Rezaei and Salomons, 2022). Raveendhran and Fast (2021) studied computerized tracking tools and show that they undermined workers' motivation when perceived as controlling. In our study, the evaluation of work performance by machines did not rely on information tracking. Instead, algorithms automatically computed an accurate measure of performance for each worker. Our setup purposefully discards the issue of automated evaluation tracking and its potentially negative effects. Fumagalli, Rezaei and Salomons (2022) developed a novel design to evaluate workers' preferences for human and algorithmic recruiters. They used an online incentivized experiment in which they elicited the willingness-to-pay of workers for each type of recruiter. Notably, the authors did not use deception so that workers who stated a preference for one type of recruiter were more likely to be evaluated by their preferred recruiter. They showed that human and algorithmic recruiters were rated differently even when they used the exact same information for their hiring decisions. Furthermore, low performers preferred human recruiters whereas high performers preferred algorithmic recruiters because they were perceived as putting more weight on task performance. Unlike Fumagalli, Rezaei and Salomons (2022), we study behavioral reactions to dismissal algorithms, which are assessed using performance in an incentivized task.³

³ Recent exceptions to the dominant approach focusing on attitudes and preferences (see Mahmud et al., 2022 for a review) rather than on work performance are the lab experiment of Strobel (2021) and the field experiment of Bai et

2.1.2. Algorithm aversion

Beyond algorithmic supervision and evaluation, an extensive literature has studied the reluctance of workers to accept algorithmic decisions in numerous applications ranging from medical diagnosis to robo-advisors (see Mahmud et al., 2022 for a review). Our work contributes to this literature on algorithm aversion by considering a new application (dismissal decisions) and a new experimental paradigm. Furthermore, we study a novel psychological mechanism based on spite (see Fehr and Schmidt, 1999; Bolton and Ockenfels, 2000; Charness and Rabin, 2002; Fehr and Fischbacher, 2002) to explain the reaction of dismissed workers to algorithmic and human decisions.

2.1.3. Ethical decisions

In a recent work, Hidalgo et al., (2021) assessed the ethical dimension of decisions involving humans or algorithms using an experimental design in which they described scenarios (vignettes) to participants. Although no scenarios contemplated the type of dismissal task we envision here, many considered work-related situations. For example, several scenarios assessed how people perceived a biased recruiting system that was implemented by a human or an algorithm. The results show that people judged humans on their perceived intentions whereas they evaluated machines on final outcomes. As a result, people judged humans more negatively when their decisions impacted others negatively and when they were seen as responsible for their actions. This was, for example, the case in scenarios evoking a biased recruiting system.

Dismissal decisions are ethical decisions because they directly impact the payoffs and psychological well-being of workers (Theodossiou, 1998; McKee-Ryan et al., 2005; Paul and Moser, 2008). Furthermore, dismissal decisions involve fairness issues because they are subject to numerous biases (Haidt, 2012) and because people might disagree on which rule to apply (Konow, 2003).

2.2. Hypotheses

In our setup, fired workers will tend to attribute negative intentions to human bosses whereas they might not do so for algorithms. This can even occur in cases in which there are no biases in

al., (In Press). Strobel (2021) uses a principal-agent experiment to show that when an algorithm is endogenously selected to enforce the minimum level of transfer of the agent to the principal then transfers are lower than when a human principal is selected. Bai et al., (In Press) report that an algorithmic assignment of tasks is not only perceived as fairer than a human assignment, but it also leads to productivity gains.

dismissal decisions but fired workers perceive it otherwise. In our setup, this is likely to occur in easy and measurable tasks. Given that workers perform equally well in these easy tasks and performance is the only available information, human bosses and algorithms will dismiss a worker at random. Yet, dismissed workers will likely perceive the outcome of the decision as unfair because they achieve the same level of performance as the other worker but obtain a lower pay (see Konow, 2003). This dismissal procedure contrasts with the basic principles of distributive justice (Adams, 1965; Nozick, 1974; Konow, 2003). It is not consistent with *equality* because workers are not paid the same amount regardless of performance.⁴ It also contrasts with *equity* because pay is not proportional to work effort, and it cannot be explained by *efficiency* principles because bosses cannot identify which worker has the highest ability level.

Workers tend to perceive the unfair dismissal as being intentional in the case of the human boss but not in the case of the algorithm (De Cremer, 2020; Hidalgo et al., 2021). Because people perceive a disadvantageous distribution of income more negatively when it is intentional, they will tend to react more negatively to human than algorithmic decisions (Blount, 1995; Kagel, Kim and Moser, 1996; Falk, Fehr and Fischbacher, 2008). As a result, workers who are dismissed by a human boss will immediately lower their effort in line with negative reciprocity motives (Fehr, Gächter and Kirchsteiger, 1997; Bolton and Ockenfels, 2000; Fehr and Gächter, 2000; Orhun, 2018).

Furthermore, a series of works have shown that reciprocity is magnified by emotional arousal (Ben-Shakhar et al., 2007; Hopfensitz and Reuben, 2009; Bolle, Tan and Zizzo, 2014; Dickinson and Masclet, 2015). The emotional reaction associated with reciprocal behavior is likely to be more pronounced when workers interact with human bosses than with algorithms because perceived intentions trigger a larger emotional response (Lee, 2018).⁵

In sum, we expect fired workers to react more negatively to human than to algorithmic dismissals. By contrast, workers who are not dismissed will tend to react more positively to human decisions

⁴ The equality principle is relevant to our context in which limited information (i.e., performance on the task) is available to the parties. In richer contexts, need (a concern for the well-being of the least well-off members of society) rather than equality is likely to be a more prevalent justice principle (Konow, 2003).

⁵ At the brain level, specific brain areas related to mentalizing (temporoparietal junction and medial prefrontal cortex), and social motivation (hypothalamus and amygdala) are only involved in reciprocal interactions with other humans (Rauchbauer et al., 2019).

because they tend to trigger positive reciprocity motives. This asymmetry of reaction between fired and non-fired workers is captured in Hypothesis *li*.

In Hypothesis *lii*, we consider the case of hard and measurable tasks that are characterized by a high level of heterogeneity in performance across workers. This is a case in which both human and algorithmic dismissal decisions hinge upon observed performance differences across workers. That is, dismissals are consistent with equity and efficiency principles that people tend to agree upon (Hidalgo, 2021).⁶ It follows that workers perceive human dismissals as fairer than in the case of the easy and measurable task in which dismissal decisions are arbitrary and no distributive justice principles apply (Konow, 2003).⁷ We thus expect the difference in workers' reaction between human and algorithmic dismissals will be less pronounced in the hard and measurable task than in the easy and measurable task.

In Hypothesis *liii*, we consider non-measurable tasks for which human bosses cannot accurately evaluate workers' performance. It implies that human bosses, unlike algorithms, will sometimes dismiss the most productive workers, thus contradicting equity and efficiency principles. This is a case in which machine evaluations can be preferred to human ones because they use more reliable information on a well-accepted performance metric (see Dietvorst, Simmons and Massey, 2015; Cowgill, 2018; Cowgill and Tucker, 2020; Kleinberg et al., 2018, 2020; Fumagalli, Rezaei and Salomons, 2022). Unlike the easy and measurable task, human bosses and algorithms will make different decisions in the non-measurable tasks because they process information differently. We thus expect fired workers will perceive human decisions as less fair than those of algorithms leading them to reduce their level of effort on the task. This conjecture is also in line with the fact that people are more likely to use algorithms when a task is more difficult (Bogert, Schechter, Watson, 2021; Prahl and Van Swol, 2021) and when it becomes clear that algorithms are more accurate than humans (Dietvorst, Simmons and Massey, 2015; Filiz et al., 2021).

⁶ Workers could still oppose these principles and favor equality. However, equality cannot be obtained in our design because one of the two workers must be dismissed and pay inequality will thus always be observed in the second stage. In addition, workers never produced the same value on the hard and measurable task, thus not providing a natural justification for the equality principle.

⁷ In our experiment, the dismissal decisions in the hard and measurable task were indeed associated with a higher perception of distributive justice than those in the easy and measurable task, regardless of the treatment (Rank Sum Tests, p -values < 0.001)

We state Hypothesis 1 as follows. For clarity purposes, we refer to the difference between the *Algo* and the *Human* treatments for a given dependent variable (production, perceived distributive justice or work satisfaction) as the *Algo treatment effect*.

Hypothesis 1 (Human vs algo dismissals and task measurability)

i) For measurable tasks, the *Algo treatment effect* of fired workers will be positive for all dependent variables. The difference in the *Algo treatment effect* between fired and non-fired workers will also be positive for all dependent variables.

ii) The effects in *i)* will be less pronounced for hard tasks that are characterized by a high heterogeneity in task performance.

iii) For non-measurable tasks, the *Algo treatment effect* of fired workers will be positive for all dependent variables. The difference in the *Algo treatment effect* between fired and non-fired workers will also be positive for all dependent variables.

In Hypothesis 2, we consider the case of tasks for which dismissal decisions are likely to be perceived as biased, that is based on dimensions that are not directly relevant for performing the job (Fumagalli, Rezaei and Salomons, 2022). In our experiment, we design bias tasks that are easy enough so that all workers can complete them perfectly so that performance metrics cannot drive decisions. Yet, workers must reproduce different patterns. Half the workers are asked to reproduce a pattern showing the word ‘FEMME’ (i.e., Woman) while the others reproduced the word ‘HOMME’ (i.e., Man).

Although both types of workers perform equally well on the task, being fired by a human can be perceived as particularly unfair because of two reasons.⁸ The first reason is similar to the one evoked for the easy and measurable task. Fired workers can see dismissals as unfair because they perform the same as the other worker. As a result, dismissal decision conflict with the standard distributive justice principles of equality, equity and efficiency (Konow, 2003). The second reason is that fired workers can perceive as unfair being assigned to the ‘Man’ (‘Woman’) pattern when the human boss favors the other pattern. The unfair decision of the human boss is likely to be perceived as intentional unlike the decision of the algorithmic boss (Hidalgo et al., 2021). Furthermore, fired workers are more likely to attribute intentions to the decision of human bosses

⁸ Note that, generally, human decisions tend to be perceived as more biased than algorithmic ones (Newman, Fast and Harmon, 2020).

in bias tasks than in the easy and measurable task because there exists a salient feature (gender) that is inherent to the task and which is often identified as a source of discrimination in work pay and hiring policies (Weichselbaumer and Winter-Ebmer, 2005; Blau and Kahn, 2017 and Lambrecht and Tucker, 2019). In sum, we expect the difference in fired workers' reaction between human and algorithmic dismissals to be especially pronounced in bias tasks.

Hypothesis 2 (Human vs algo dismissals in bias tasks)

For bias tasks, the *Algo treatment effect* of fired workers will be positive for all dependent variables. The difference in the *Algo treatment effect* between fired and non-fired workers will also be positive for all dependent variables.

In Hypothesis 3, we consider the case in which one worker is assigned to a hard task (handicap task) that is impossible to complete perfectly while the other worker is assigned to an easy task that can be completed perfectly by any participant. As a result, workers who are assigned to the hard task always produce less than those assigned to the easy task. Human bosses thus have a material interest to dismiss the handicapped workers who produce less and generate less revenues. However, human bosses can also decide to discard the efficiency principle and follow equity, thus valuing the higher level of effort exerted by the handicapped worker (see e.g., Leventhal and Michaels, 1971 and Ruffle, 1998). Unlike human bosses, our dismissal algorithms always demote the least productive worker.

In the case of handicap tasks and similarly to bias tasks, we expect dismissed workers to react more negatively to human than to algorithmic dismissals (see Hypothesis 3). This is the case because bosses and workers are likely to disagree on which distributive principle should apply (equity or efficiency). Humans tend to develop self-serving fairness norms in situations in which several principles (here efficiency and equity) conflict, thus perceiving as fairer those rules that are beneficial to them (Knez and Camerer, 1995; Kagel, Kim and Moser, 1996; Konow, 2000; Gächter, and Riedl, 2005; Bolton and Ockenfels, 2008; Bejarano, Corgnet and Gómez-Miñambres, 2021). If bosses were driven by self-interest, they would tend to follow efficiency principles and base their dismissal decisions on output. In that case, dismissed workers would tend to think that equity principles should apply so that effort levels are taken into account. By contrast, if bosses followed equity principles, then dismissed workers who produced a high level of output would likely

embrace efficiency principles. As a result, regardless of the boss decision, fired workers will tend to perceive the decision as unfair. We state Hypothesis 3 as follows.

Hypothesis 3 (Human vs algo firing dismissals in handicap tasks)

For handicap tasks, the *Algo treatment effect* of fired workers will be positive for all dependent variables. The difference in the *Algo treatment effect* between fired and non-fired workers will also be positive for all dependent variables.

Finally, we expect our results to be impacted by workers' general attitude towards robots as measured using the Negative Attitudes Towards Robots scales (NARS, henceforth) (Syrdal et al., 2009). Although Dietvorst, Simmons and Massey (2015) interpret people's willingness to avoid algorithms whenever they make observable mistakes as some form of algorithm aversion (see Mahmud et al., 2022 for a review), they have not assessed the moderating role of general attitudes towards robots. We are not aware of any study, regardless of the methodology, assessing the moderating role of NARS on workers' reaction to dismissal algorithms. We believe it is a relevant question with practical implications regarding the implementation of algorithmic dismissals. If people who score high on NARS are reluctant to be dismissed by algorithms, this should be taken into account by practitioners who might first want to experiment with workers who have positive attitudes towards robots and in industries that are generally favorable to the use of robots (Wang et al., 2010). As a result, we conjecture in Hypothesis 4 that the positive reaction to algorithmic compared to human dismissals will be even more pronounced for fired workers who have positive attitudes towards robots (i.e., a low NARS score).⁹

Hypothesis 4 (Human vs algo firing: Algorithm aversion)

For all tasks, the positive *Algo treatment effect* will be more pronounced for fired workers who have a low NARS score. The positive difference in the *Algo treatment effect* between fired and non-fired workers will be more pronounced for those with a low NARS score.

⁹ We pre-registered the hypotheses in AsPredicted (# 96732). The document is available here: <https://aspredicted.org/zj4dd.pdf>. For clarity, in our hypotheses section, we use the terms work satisfaction, distributive justice and handicap tasks instead of intrinsic motivation, procedural justice and discrimination tasks in the pre-registered document. We also modified the ordering of the sentences and defined the *Algo treatment effect* to simplify the statement of the hypotheses. These are cosmetic changes that do not affect the content of the scales, the nature of the tasks and the tests used. For the sake of parsimony, we also decided to state and test the second part of the pre-registered hypothesis (Hypothesis 4ii) in a separate Appendix C.

3. Design

A crucial feature of our design is that we did not use deception. We implemented an actual algorithm that, unlike the popular ‘Wizard of Oz’ approach (Cross and Ramsey, 2021), was not under the control of a human monitor during the experiment. The absence of deception is standard in the economics literature (e.g., Fumagalli, Rezaei and Salomons, 2022) but not in the human-robot interaction literature (see Riek, 2012). We implemented an actual algorithm in a physical lab facility that is known to forbid deception to make sure participants believed our experimental manipulations (Jamison, Karlan and Schechter, 2008).¹⁰ Our design builds on the ‘computer-as-player’ approach used in economic experiments in which one or several human participants in the baseline treatment are replaced by a machine player in the automated treatment (see Chugunova and Sele, 2020; March, 2021 for reviews). Thus far, this literature has focused on the impact of automated players on teamwork and cooperation (Crandall et al., 2018; Corgnet, Hernán-González and Mateo, 2019; Alekseev, 2020; Traeger et al., 2020), trust (Schniter, Shields and Sznycer, 2020), and moral decisions impacting other participants’ payoffs (Gogoll and Uhl, 2018; Kirchkamp and Stroebel, 2019; Strobel, 2021).

3.1. The Dismissal Game (DSG)

We designed a new experimental paradigm to study dismissal decisions in a laboratory experiment that we will refer to as the Dismissal Game (DSG, henceforth). The DSG was composed of two stages. In the first stage, two workers completed a real-effort task for a 5-euro fixed pay. By completing the task, workers produced monetary value that was directly transferred to a boss.

At the end of the first stage, the boss decided, after observing feedback on workers’ performance, which one to dismiss and which one to keep. The dismissed worker was paid a fixed salary of 1 euro instead of 5 in the second stage of the game whereas the other worker received the same pay as in the first stage (5 euros). Both workers completed the exact same task in the second stage as in the first stage. Assessing the change in workers’ performance between the two stages of the game allows us to study the behavioral consequences of dismissals. In our game, the second stage can be seen as the notice period which is the time that elapses between the announcement of a dismissal decision and its actual implementation. This period typically ranges from several weeks

¹⁰ On a side note, we designed this protocol to collect data in March 2020. Because of pandemic restrictions we had to wait for two years to conduct it in the laboratory. We did not want to collect data online to make sure participants believed they were interacting with other human participants in the *Human* treatment.

in the US and Canada to several months in most European countries (e.g., Venn, 2009). More fundamentally, the second stage of the game aims at capturing potential costs associated with a negative behavioral reaction of workers to dismissals that can impact employers' revenues. These costs could, for example, be related to legal procedures triggered by resentful workers who perceived they have been unfairly dismissed (Lind et al. 2000; Goldman, 2003). In sum, workers' production in Stage 2 offers us a behavioral proxy of organizational justice (Colquitt et al., 2013). Regarding the fired worker, the slash in compensation from \$5 to \$1 aims at capturing the monetary costs of being dismissed such as the search costs associated with finding a new job. Our design is equivalent to one that continues to pay the dismissed worker \$5 in the second stage and inflicts a \$4 penalty to the worker at the end of the period.

The two-stage game was repeated across seven periods. Participants kept the same role (boss or worker) for the whole duration of the experiment. However, they were randomly matched with other participants in each period.¹¹ In each period, workers were assigned to a different task.

3.2. Tasks

Before starting the first period and before being assigned a role, all participants completed an ability task (see Appendix A.1) that had a sufficiently high level of difficulty so that no participant could achieve the maximum level of performance, and this was indeed the case. As for the tasks used in the DSG, the 90-second ability task required reproducing a colored pattern represented in a 20×20 grid by coloring pixels (cells) using five possible colors (blue, grey, green, yellow and black).¹² A cell filled with the correct color produced a value of 2¢ whereas a cell filled with the wrong color implied a penalty of -1¢. The average (standard deviation) pay for the ability task was 54.86¢ (27.85¢). Negative pay was not allowed.¹³

Unlike the ability task, the tasks for the DSG could differ across workers. Across periods, tasks also differed in their level of difficulty and in the measurability of performance. Workers produced monetary value for the boss whenever they completed these tasks, regardless of the period or stage of the DSG. This value was calculated as for the ability task so that each cell filled with the correct (wrong) color produced 2¢ (-1¢) for the boss. Producing negative value was not permitted. Before

¹¹ This is not a perfect stranger matching design but, given a average session size of 21, the likelihood of being matched with the same employer and employee (or the same employees) in two consecutive periods is only 2%. The likelihood that you will never encounter the same pairs in any of the following periods is about 60%.

¹² Colors were chosen to limit common color-blindness issues.

¹³ This constraint was not binding because all workers produced a strictly positive amount.

making their dismissal decision, bosses observed the output produced by each worker for 25 seconds so that they could estimate the value of the tasks. Each task lasted 90 seconds.

Measurable and non-measurable tasks

The task completed in the first period (see Task 1 in Appendix A.1) was the same for both workers and was both easy and measurable. This task was easy because we expected all participants to complete it perfectly, and this was actually the case in the first stage. Like Task 1, the performance on Task 3 was also easily measurable but the task was hard to complete entirely. Indeed, no workers filled the 400 cells correctly. Unlike Task 1, we expected a high level of heterogeneity in workers' performance in Task 3, and this was the case. Like Task 3, Task 6 was hard, and no workers completed the whole pattern perfectly. It differed from Task 3 because measuring the exact performance on the task was extremely difficult due to the complexity of the color pattern. We refer to Task 6 as non-measurable.

Bias tasks

Tasks 2 and 5 were easy because we expected most workers to complete them perfectly, which indeed occurred in 94.0% of the cases in the first stage. These tasks differed from Tasks 1, 3 and 6 because each of the two workers had to complete a slightly different task (see Appendix A.1). One of the workers had to reproduce the word 'Homme' (i.e., 'Man') whereas the other one reproduced the word 'Femme' (i.e., 'Woman'). We added cells at the bottom of the two patterns so that each task required filling the same number of cells (60), and thus led to the same value for the boss when duly completed. The 'Woman' and 'Man' tasks were thus of the same difficulty (Rank Sum Test, p -value = 0.226 for a test comparing performance across tasks in Stage 1). We refer to Tasks 2 and 5 as bias tasks because dismissal decisions could be perceived as impacted by gender preferences that are unrelated to the performance of each worker. Because these tasks were designed so that most workers completed them perfectly, the gender word was the only distinctive feature available to make dismissal decisions.

Handicap tasks

As for bias tasks, handicap tasks (Tasks 4 and 7) required each worker to complete a different task. One of the two tasks was easy, and the other was hard. The easy task was similar to Task 1 and all participants were expected to complete it perfectly, which indeed happened in 93.0% of the cases in the first stage. The hard task was such that no worker could complete it perfectly, which was

indeed always the case. As a result, the worker endowed with the easy task always produced more than the one endowed with the hard task in Stage 1. Because the tasks faced by workers in Stage 2 were the same as in Stage 1, bosses knew that workers endowed with the easy task would tend to produce more value than those endowed with the hard task. This was indeed the case in all instances.¹⁴

3.3. Treatments

In addition to the within-variation in the nature of the task across periods, we considered two between-subject treatments. In the *Human* treatment, the dismissal decision at the end of Stage 1 was made by the human boss. Instructions stated explicitly that the boss was another human participant selected at random (see Appendix A.2).

In the *Algo* treatment, an algorithm automatically dismissed the worker who had produced the least valuable drawing in the first period. This automated rule was solely based on workers' performance and did not account for workers' effort on the task. This allowed us to implement a simple rule that could be easily communicated and understood by participants to avoid any negative reaction to automation due to lack of transparency (Shin and Park, 2019; De Cremer, 2020). The automated rule was known to workers and stated as follows in the instructions:

“At the end of the first period, a pre-programmed robot will decide which of the two workers P1 or P2 will be kept at his workstation for the second period. (...) To make its decision, the robot will compute the value of the pattern produced by each worker. The worker who produces the most value will be maintained, and the other worker will be dismissed. In case of a tie, the robot selects a worker at random.” In the *Algo* treatment, we kept the number of human participants constant so that, as in the *Human* treatment, one of the participants was assigned the role of Participant *B* (boss) and received the monetary value that workers produced on the task. This passive player allowed us to control for the effect of social preferences that are unrelated to automation. For example, in the *Human* treatment, workers could attempt to reduce inequality in pay across participants by exerting effort to boost the revenues of the boss.

¹⁴ This observation is evidently altered by endogeneity issues because Stage 2 performance is likely to be affected by dismissal decisions, which, in turn, are likely affected by the difficulty of the task a worker was endowed with.

3.4. Procedures

Our design was conducted in a major lab in France. We pre-registered the design, hypotheses, and analyses in AsPredicted (# 96732). In line with our pre-registered design, we recruited 237 participants (39 triplets for the *Human* treatment and 40 for the *Algo* treatment) from a subject pool of more than 2,000 students across all disciplines. We conducted a total of 12 sessions with either 15, 18, 21, 24 or 27 participants. After the round of instructions and before the experiment started, we asked participants to complete a 7-question comprehension quiz (see A.2 in Appendix A). The quiz was incentivized so that participants who answered all questions correctly in their first attempt received a 1-euro bonus. Overall, participants answered 76.2% of the questions correctly in their first attempt. At the end of the experiment, we conducted a 5-minute questionnaire eliciting social preferences (see Bartling et al., 2009 and Corgnet, Espín and Hernán-González, 2015), the Negative Attitudes Towards Robots scales (NARS, henceforth) (Syrdal et al., 2009) and basic demographics (see A.2 in Appendix A).

Regarding payments for the main phase of the experiment (DSG), one of the 7 periods was selected at random. Regardless of the selected period, a worker who was maintained (dismissed) in that period received 10 (6) euros. On average, bosses earned 8.6 euros in the DSG. Final payments, including the earnings on the ability task, the social preferences test, the bonus for the comprehension quiz and a show-up fee of 5 euros, were on average equal to 15.0 euros for an experiment lasting one hour.

4. Results

4.1. Dependent variables

Following our pre-registration plan, we considered three dependent variables.¹⁵ Our behavioral variable is the difference in workers' production (in cents) between stages. We also consider two measures of attitudes that were elicited at the end of each period.¹⁶ Distributive justice scores were based on a 3-item pay scale adapted from Colquitt (2001) (Cronbach $\alpha = 0.925$). Our 4-item scale

¹⁵ In our pre-registration, we contemplated using both non-parametric tests and panel regressions to test our hypotheses. Because of the dynamic structure of our data, panel analyses are more appropriate, and we thus report these results here. In the cases in which we had only one observation per worker, we report OLS estimates (see Tables 2, B2 and B3). Similar findings are obtained using non-parametric tests and are not reported here for the sake of parsimony.

¹⁶ Because our focus was on behavior, attitudes were never measured before work production so that they would not impact work behavior.

of work satisfaction was based on the intrinsic motivation inventory (see Ryan, 1982) (Cronbach $\alpha = 0.911$). The full scales are shown in Appendix A.2.

Following our hypotheses, for each dependent variable, we focus on the *Algo treatment effect*, which captures the effect of algorithmic dismissals. The difference in the *Algo treatment effect* between fired and non-fired workers is captured by the interaction term ‘Algo \times Fired’, where ‘Algo’ (‘Fired’) is a dummy variable that takes value one for the *Algo* treatment (for fired workers) and value zero otherwise. In line with our hypotheses, we also assess the sign and significance of the *Algo treatment effect* of fired workers by checking whether the following combination of coefficients is equal to zero: ‘Algo \times Fired’ + ‘Algo’. According to our hypotheses, the coefficients for ‘Algo \times Fired’ and for the combination ‘Algo \times Fired’ + ‘Algo’ should be positive and significant. For the sake of completeness, we also report the *Algo treatment effect* of non-fired workers by checking the coefficient for ‘Algo’ as well as the impact of being fired in the *Human* (see ‘Fired’ coefficient) and *Algo* treatments (see ‘Algo \times Fired’ + ‘Fired’). Although the main text of the results section report on all the coefficients, we made sure that regression tables were self-contained so that the reader will lose little substance by focusing on regression tables.

4.2. Aggregate results

4.2.1. Change in production

We start by presenting the aggregate results across all tasks. We first compare workers’ production in the first stage in the *Algo* and *Human* treatments to make sure that the two treatments did not significantly differ in terms of incentives effects. Reassuringly, we find no differences across treatments in Stage 1.¹⁷ This is shown in regression [1] in Table B1 (see Appendix B) where the ‘Algo’ variable is not significant in explaining production.

Overall, fired workers decreased their production more than those who were not fired in the *Human* treatment whereas no difference was observed in the *Algo* treatment (see Figure 1). Furthermore, fired workers decreased their production by a larger amount in the *Human* (Mean = -44.72ϕ , SD = 110.86ϕ) than in the *Algo* treatment (Mean = -17.47ϕ , SD = 70.84ϕ). For workers who were not fired, differences across treatments were less pronounced (Mean = 0.49ϕ , SD = 40.51ϕ in the *Human* treatment, and Mean = -4.48ϕ , SD = 52.73ϕ in the *Algo* treatment). Considering all workers,

¹⁷ Note that, unlike Raveendhran and Fast (2021), we did not implement tracking tools that could have impacted workers’ motivation differently in the *Human* and *Algo* treatments.

we find that the decrease in production was larger in the *Human* (Mean = -22.11ϕ , SD = 86.40ϕ) than in the *Algo* treatment (Mean = -10.97ϕ , SD = 62.73ϕ).

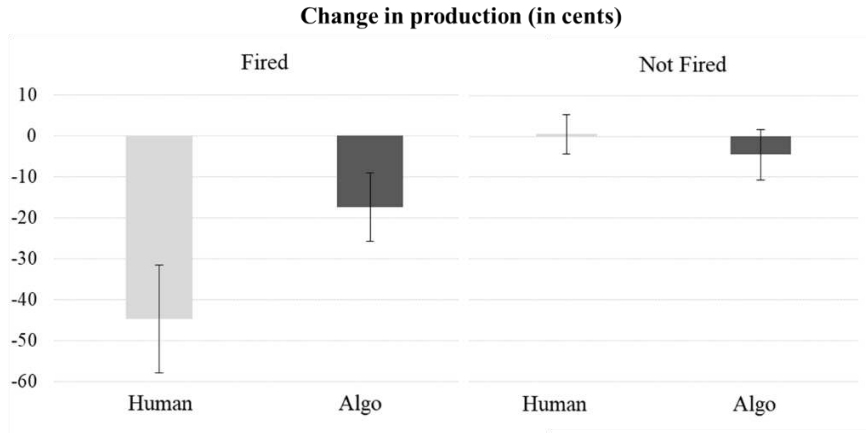


Figure 1. Average change in production between stages (in cents) across treatments (Human and Algo) for fired workers (left panel) and non-fired workers (right panel) along with 95% confidence intervals.

We test the statistical significance of these differences using panel regressions that account for the fact that we have several observations per worker. In Table 1 (regression [1]), we report that the coefficient for ‘Algo \times Fired’ is positive and significant so that algorithmic dismissals induced a more positive change in the performance of fired workers compared to non-fired workers (see Table 1, p -value = 0.001). The *Algo* treatment impacts the change in production negatively but not significantly so for non-fired workers (see negative and significant ‘Algo’, p -value = 0.191). The *Algo* treatment has a positive and significant effect on the change in production for fired workers (‘Algo \times Fired’ + ‘Algo’ = 0, p -value = 0.011). We note that the total effect of algorithm dismissals on the change in production, which can be estimated as $\frac{1}{2}$ (‘Algo \times Fired’) + ‘Algo’, is positive yet not significant (p -value = 0.117). Finally, the decrease in production for being fired is significant in the *Human* treatment (see ‘Fired’, p -value < 0.001) but fails to reach significance at standard levels in the *Algo* treatment (see ‘Algo \times Fired’ + ‘Fired’ = 0, p -value = 0.052).

Table 1. Change in production, perceived distributive justice, and work satisfaction.

This table presents the results of linear panel regressions with random effects and clustered standard errors (in parentheses) at the session level using wild bootstrapping techniques as in Cameron and Miller (2015) with 5,000 replications, which is recommended given our number of clusters (12). The *Algo* (*Fired*) Dummy takes value one for the *Algo* treatment (when a worker has been fired) and value zero otherwise. Ability is the performance on the task completed by all participants at the beginning of the experiment. We include fixed effects for task types.

Dependent Variable	Change in production [1]	Perceived distributive justice [2]	Work satisfaction [3]
Constant	27.352** (11.119)	16.324*** (0.509)	19.686*** (0.848)
<i>Algo</i> Dummy	-6.857 (5.336)	-1.126*** (0.374)	0.186 (0.611)
<i>Fired</i> Dummy	-42.728*** (9.831)	-5.377*** (0.370)	-2.743*** (0.276)
<i>Algo</i> × <i>Fired</i> Dummy	33.976*** (10.719)	1.873*** (0.537)	1.287*** (0.439)
Period	-5.712*** (1.366)	0.103 [◇] (0.056)	-0.242*** (0.070)
Ability	-0.054 (0.062)	-0.002 (0.005)	0.011 (0.008)
Male Dummy	-16.082*** (5.131)	-0.626 [◇] (0.363)	-1.802*** (0.620)
N	1,106	1,106	1,106
R ²	0.092	0.410	0.099
Prob > χ^2	<0.001	<0.001	<0.001
P-values (F-Tests)			
<i>Algo</i> × <i>Fired</i> + <i>Algo</i>	0.011	0.152	0.025
<i>Algo</i> × <i>Fired</i> + <i>Fired</i>	0.052	<0.001	<0.001

*** Significant at the 0.01 level; ** at the 0.05 level; [◇] for *p*-values in (0.05, 0.10).

4.2.2. Distributive justice and work satisfaction

In line with the behavioral data on performance, fired workers reported lower distributive justice (Mean = 9.67, SD = 3.53) and work satisfaction (Mean = 16.81, SD = 5.82) compared to non-fired workers (Mean = 14.19, SD = 3.00; Mean = 19.18, SD = 5.54, respectively). These differences were larger for the *Human* treatment (Mean = 9.25, SD = 3.33 vs Mean = 14.77, SD = 2.59 for distributive justice, Mean = 15.92, SD = 6.14 vs Mean = 19.14, SD = 5.89 for work satisfaction) than for the *Algo* treatment (Mean = 10.09, SD = 3.68 vs Mean = 13.63, SD = 3.26 for distributive justice, Mean = 17.69, SD = 5.35 vs Mean = 19.23, SD = 5.19 for work satisfaction) (see Figure 2). Furthermore, fired workers evaluated distributive justice and work satisfaction to be lower in the *Human* than in the *Algo* treatment. By contrast, workers who were not fired reported lower distributive justice in the *Algo* than in the *Human* treatment and similar levels of work satisfaction across treatments. Considering all workers, we find similar levels of distributive justice (Mean =

12.01, SD = 4.06 vs Mean = 11.86, SD = 3.90) and work satisfaction (Mean = 17.53, SD = 6.22 vs Mean = 18.46, SD = 5.32) in the *Human* and *Algo* treatments.

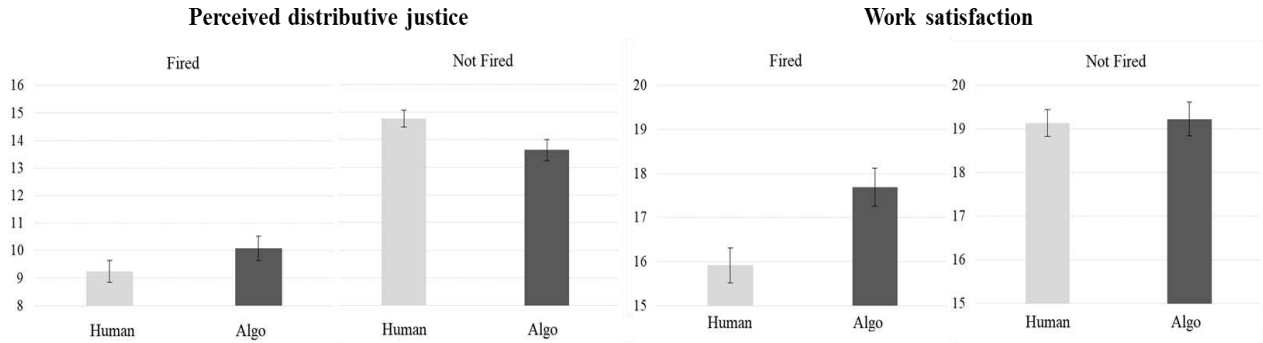


Figure 2. Average perceived distributive justice [work satisfaction] across treatments (Human and Algo) for fired workers and non-fired workers on the left [right] panel along with 95% confidence intervals.

Using panel regressions, we first report that the coefficient for ‘Algo × Fired’ is positive and significant in regression [2] of Table 1 (p -value < 0.001) so that algorithmic dismissals induce higher perceived distributive justice for fired workers than for non-fired workers. Yet, the *Algo* treatment has a non-significant effect on fired workers (‘Algo × Fired’ + ‘Algo’ = 0, p -value = 0.152) and a significantly negative effect for non-fired workers (see negative and significant ‘Algo’, p -value = 0.003). We note that the total effect of algorithm dismissals on perceived distributive justice is not significant (p -value = 0.597). Finally, perceived distributive justice is significantly lower for fired than for non-fired workers in the *Human* (see ‘Fired’, p -value < 0.001) and *Algo* treatments (see ‘Algo × Fired’ + ‘Fired’ = 0, p -value < 0.001).

Regarding work satisfaction, the coefficient for ‘Algo × Fired’ is also positive and significant (see regression [3] in Table 1, p -value = 0.003). The *Algo* treatment does not impact work satisfaction for non-fired workers (see ‘Algo’, p -value = 0.761) but it increases fired workers’ satisfaction significantly (‘Algo × Fired’ + ‘Algo’ = 0, p -value = 0.025). We note that the total effect of algorithm dismissals on work satisfaction is positive yet not significant (p -value = 0.166). Finally, work satisfaction is significantly lower for fired workers in the *Human* (see ‘Fired’, p -value < 0.001) and *Algo* treatments (see ‘Algo × Fired’ + ‘Fired’ = 0, p -value < 0.001).

4.3. Hypotheses tests

In this section, we analyze differences across task types, thus testing each of our first three hypotheses sequentially.

4.3.1. Hypothesis 1

Performance

Hypothesis Ii

In line with Hypothesis *Ii*, we report that for measurable tasks the coefficient for ‘Algo × Fired’ is positive and significant (see regression [1] in Table 2, p -value = 0.047). Yet, the positive effect of the *Algo* treatment on the change in production for fired workers is not significant (‘Algo × Fired’ + ‘Algo’ = 0, p -value = 0.151). The *Algo* treatment does not impact the change in production for non-fired workers (see ‘Algo’, p -value = 0.411). We note that the total effect of algorithm dismissals on the change in production is not significant (p -value = 0.415). Finally, the change in production is significantly lower for fired workers compared to non-fired workers in the *Human* treatment (see ‘Fired’, p -value < 0.001) but not in the *Algo* treatment (see ‘Algo × Fired’ + ‘Fired’ = 0, p -value = 0.121).

Hypothesis Iii

As expected, the easy and hard measurable tasks differed regarding the heterogeneity in workers’ performance in the first stage. In the easy measurable task, workers always produced the maximum value except in one case. In the hard measurable task, workers never produced the same value in the first stage. Furthermore, the difference in performance was strictly greater than 10¢ (which is equivalent to 10 correctly filled cells) in more than 90% of the cases. We also found that when workers’ performance differed in Stage 1 in the hard measurable task, human bosses picked the highest performer in 92.3% of the cases, which is similar to the rule applied by the algorithm (100%).

To test Hypothesis *Iii*, we consider the easy and hard measurable tasks separately in regressions [2] and [3]. In line with Hypothesis *Ii*, ‘Algo × Fired’ continues to be positive and significant for the easy measurable task (see regression [2] in Table 2, p -value = 0.048). The *Algo* treatment has a positive and significant effect on the change in production for fired workers (‘Algo × Fired’ + ‘Algo’ = 0, p -value = 0.041) but not for non-fired workers (see ‘Algo’, p -value = 0.855). We note

that the total effect of algorithm dismissals on the change in production is positive and significant (p -value = 0.045) for the easy measurable task. Finally, the change in production is negative, yet not significant, for fired workers in the *Human* treatment (see ‘Fired’, p -value = 0.059), and positive and not significant in the *Algo* treatment (see ‘Algo × Fired’ + ‘Fired’ = 0, p -value = 0.604).

In line with Hypothesis *l*_{ii}, ‘Algo × Fired’ is not significant for the hard measurable task (see regression [3] in Table 2, p -value = 0.128). Unlike the easy measurable task, the *Algo* treatment does not impact the change in production for fired workers (see ‘Algo × Fired’ + ‘Algo’ = 0, p -value = 0.260) and non-fired workers (see ‘Algo’, p -value = 0.471). We note that the total effect of algorithm dismissals on the change in production is not significant (p -value = 0.655). It is also interesting that the positive coefficient associated with ‘Algo × Fired’ is substantially decreased if we consider cases in which the difference in performance between the two workers was large enough to be visually detected by workers on their feedback screen.¹⁸ Finally, the change in production is significantly lower for fired workers in the *Human* treatment (see ‘Fired’, p -value = 0.005) but not in the *Algo* treatment (see ‘Algo × Fired’ + ‘Fired’ = 0, p -value = 0.189).

These findings provide support for Hypothesis *l*_{ii} according to which the positive effect of the *Algo* treatment on the change in production of fired workers will be observed for the easy but not for the hard measurable task.

¹⁸ If we consider the cases in which the difference between the production of the two workers was at least 10¢ (that is 10 correctly filled cells), which corresponds to more than 90% of the cases, the coefficient for ‘Algo × Fired’ was halved (from 54.61 to 26.67). The coefficient even turns out to be negative if we consider all instances in which the difference between the production of the two workers is above the median difference (i.e., 80¢).

Table 2. Change in production and task measurability.

This table presents the results of panel <OLS> regressions with random effects and clustered standard errors (in parentheses) at the session level using wild bootstrapping techniques as in Cameron and Miller (2015) with 5,000 replications, which is recommended given our number of clusters (12) for regression [1] <[2], [3] and [4]>. The Algo (Fired) Dummy takes value one for the *Algo* treatment (when a worker has been fired) and value zero otherwise. Ability is the performance on the task completed by all participants at the beginning of the experiment. We include fixed effects for task type in regression [1].

Dependent Variable	Change in production			
	Measurable tasks [1]	Easy & Measurable task [2]	Hard & Measurable task [3]	Non-measurable task [4]
Constant	6.876 (14.537)	3.984 (3.590)	19.372 (29.596)	11.271 (6.259)
Algo Dummy	-9.171 (11.024)	0.228 (1.245)	-17.102 (23.747)	-6.264 (6.309)
Fired Dummy	-51.808*** (18.358)	-21.628 ^o (11.453)	-79.612*** (28.257)	-22.752** (9.344)
Algo × Fired Dummy	40.454** (20.118)	23.641** (11.936)	54.611 (35.886)	25.005** (11.926)
Ability	0.016 (0.180)	-0.099 (0.061)	0.110 (0.364)	-0.071 (0.105)
Male Dummy	-7.222 (8.853)	4.208 (3.036)	-16.135 (17.201)	-12.168** (5.490)
N	316	158	158	158
R ²	0.049	0.101	0.063	0.067
Prob > χ^2	0.006	0.165	0.026	<0.001
<u>P-values (F-Tests)</u>				
Algo × Fired + Algo	0.151	0.041	0.260	0.170
Algo × Fired + Fired	0.121	0.604	0.189	0.748

*** Significant at the 0.01 level; ** at the 0.05 level; ^o for *p*-values in (0.05, 0.10).

Hypothesis Iiii

In the case of the non-measurable task, as expected, human bosses dismissed the worst performer only in 70.3% of the cases, much less often than algorithms (100%) and much less often than human bosses in the hard measurable task (92.3%, Proportion Test, *p*-values < 0.001).¹⁹ Yet, they dismissed the least productive worker more often than chance (Proportion Test, *p*-value < 0.001).

¹⁹ We cannot be certain that human dismissal decisions are solely based on performance and so it could be that a proportion of the decisions that are consistent with workers' relative performance levels could be due to mistakes or a voluntary attempt to randomize the decision. However, the results on the hard measurable task suggest humans are willing to use relative performance measures to make their dismissal decisions when the two workers face the exact same task.

In line with Hypothesis *liii*, we find that for the non-measurable task ‘Algo × Fired’ is positive and significant (see regression [4] in Table 2, p -value = 0.036). However, the *Algo* treatment does not significantly increase the change in production for fired workers (see ‘Algo × Fired’ + ‘Algo’ = 0, p -value = 0.170). The *Algo* treatment does not significantly impact the change in production for non-fired workers either (see ‘Algo’, p -value = 0.321). We note that the total effect of algorithm dismissals on the change in production is not significant (p -value = 0.485). Finally, the change in production is significantly lower for fired than for non-fired workers in the *Human* treatment (see ‘Fired’, p -value = 0.015) but not in the *Algo* treatment (see ‘Algo × Fired’ + ‘Fired’ = 0, p -value = 0.748).

Attitudes

Hypothesis li

In line with Hypothesis *li*, we report that ‘Algo × Fired’ is positive and significant (see regression [1] in Table B2, p -value < 0.001) for measurable tasks. The *Algo* treatment increased perceived distributive justice for fired workers, yet not significantly so (see ‘Algo × Fired’ + ‘Algo’ = 0, p -value = 0.092) and decreased it for non-fired workers (see ‘Algo’, p -value = 0.035). We note that the total effect of algorithm dismissals on perceived distributive justice is not significant (p -value = 0.999). Finally, perceived distributive justice was significantly lower for fired workers in the *Human* (see ‘Fired’, p -value < 0.001) and *Algo* treatments (see ‘Algo × Fired’ + ‘Fired’ = 0, p -value < 0.001).

Regarding work satisfaction, we report that ‘Algo × Fired’ is not significant (see regression [1] in Table B3, p -value = 0.734) for measurable tasks. The *Algo* treatment increased work satisfaction for fired (see ‘Algo × Fired’ + ‘Algo’ = 0, p -value = 0.128) and non-fired workers (see ‘Algo’, p -value = 0.054) but not significantly so. We note that the total effect of algorithm dismissals on work satisfaction is positive and significant (p -value = 0.046). Finally, work satisfaction was significantly lower for fired workers in the *Human* treatment (see ‘Fired’, p -value = 0.035) but not significantly so in the *Algo* treatment (see ‘Algo × Fired’ + ‘Fired’ = 0, p -value = 0.071).

Hypothesis lii

In line with Hypothesis *lii*, we report that for the easy measurable task ‘Algo × Fired’ is positive and significant for perceived distributive justice (see regression [2] in Table B2, p -value = 0.047). The *Algo* treatment reduced perceived distributive justice for non-fired workers (see ‘Algo’, p -

value = 0.082) and increased it for fired workers (see ‘Algo × Fired’ + ‘Algo’ = 0, p -value = 0.370) but none of the effects are significant. We note that the total effect of algorithm dismissals on perceived distributive justice is not significant (p -value = 0.681). Finally, perceived distributive justice was significantly lower for fired workers in the *Human* (see ‘Fired’, p -value < 0.001) and *Algo* treatments (see ‘Algo × Fired’ + ‘Fired’ = 0, p -value < 0.001).

In contrast with Hypothesis *l*ii, we report that for the hard measurable task ‘Algo × Fired’ is positive and significant for perceived distributive justice (see regression [3] in Table B2, p -value = 0.001), as in the easy measurable task. The *Algo* treatment increased perceived distributive justice for fired workers (see ‘Algo × Fired’ + ‘Algo’ = 0, p -value = 0.066) and decreased it for non-fired workers (see ‘Algo’, p -value = 0.093) but none of these effects are significant. We note that the total effect of algorithm dismissals on work satisfaction is not significant (p -value = 0.747). Finally, perceived distributive justice was significantly lower for fired workers in the *Human* (see ‘Fired’, p -value < 0.001) and *Algo* treatments (see ‘Algo × Fired’ + ‘Fired’ = 0, p -value < 0.001).

Regarding work satisfaction, we report that for the easy measurable task ‘Algo × Fired’ is not significant (see regression [2] in Table B3, p -value = 0.968). The *Algo* treatment increased work satisfaction for non-fired (see ‘Algo’, p -value = 0.074) and fired workers (see ‘Algo × Fired’ + ‘Algo’ = 0, p -value = 0.284) but not significantly so. We note that the total effect of algorithm dismissals on work satisfaction is positive yet not significant (p -value = 0.083). Finally, work satisfaction was significantly lower for fired workers in the *Human* (see ‘Fired’, p -value < 0.001) and *Algo* treatments (see ‘Algo × Fired’ + ‘Fired’ = 0, p -value = 0.001).

For the hard measurable task ‘Algo × Fired’ is not significant (see regression [3] in Table B3, p -value = 0.386). The *Algo* treatment increased work satisfaction for non-fired workers (see ‘Algo’, p -value = 0.648) and fired workers (see ‘Algo × Fired’ + ‘Algo’ = 0, p -value = 0.139) although not significantly so. We note that the total effect of algorithm dismissals on work satisfaction is positive yet not significant (p -value = 0.051). Finally, work satisfaction was not significantly different between fired and non-fired workers in the *Human* (see ‘Fired’, p -value = 0.301) and *Algo* treatments (see ‘Algo × Fired’ + ‘Fired’ = 0, p -value = 0.915).

Hypothesis Iiii

In line with Hypothesis *Iiii*, we report that for the non-measurable task ‘Algo × Fired’ is positive yet not significant (see regression [4] in Table B2, p -value = 0.142). The *Algo* treatment reduced perceived distributive justice for non-fired workers (see ‘Algo’, p -value = 0.155) and increased it for fired workers (see ‘Algo × Fired’ + ‘Algo’ = 0, p -value = 0.568) but none of the effects are significant. We note that the total effect of algorithm dismissals on perceived distributive justice is not significant (p -value = 0.833). Finally, perceived distributive justice was significantly lower for fired workers in the *Human* (see ‘Fired’, p -value < 0.001) and *Algo* treatments (see ‘Algo × Fired’ + ‘Fired’ = 0, p -value < 0.001).

Regarding work satisfaction, we report that ‘Algo × Fired’ is positive and significant (see regression [4] in Table B3, p -value = 0.037). The *Algo* treatment increased work satisfaction significantly for fired workers (see ‘Algo × Fired’ + ‘Algo’ = 0, p -value = 0.020) but not for non-fired workers (see ‘Algo’, p -value = 0.441). We note that the total effect of algorithm dismissals on work satisfaction is not significant (p -value = 0.094). Finally, work satisfaction was significantly lower for fired workers in the *Human* treatment (see ‘Fired’, p -value < 0.001) but not in the *Algo* treatment (see ‘Algo × Fired’ + ‘Fired’ = 0, p -value = 0.731).

4.3.2. Hypothesis 2

Performance

For bias tasks, in most of the cases it was not possible for human bosses to dismiss the least productive worker because, as expected, workers produced the maximum value in most cases (93.4%) and produced the same value in 88.0% of the cases.

In line with Hypothesis 2, we find that ‘Algo × Fired’ is positive and significant (p -value < 0.001, see regression [1] in Table 3). We also find that the *Algo* treatment impacted change in production positively for fired workers (‘Algo × Fired’ + ‘Algo’ = 0, p -value = 0.012) but negatively for non-fired workers (see ‘Algo’, p -value = 0.021). We note that the total effect of algorithm dismissals on the change in production is not-significant (p -value = 0.121). Finally, the change in production is significantly lower for fired workers for the *Human* (see ‘Fired’, p -value < 0.001) and *Algo* treatments (see ‘Algo × Fired’ + ‘Fired’ = 0, p -value < 0.001).

Table 3. Change in production, bias and handicap tasks.

This table presents the results of panel regressions with random effects and clustered standard errors (in parentheses) at the session level using wild bootstrapping techniques as in Cameron and Miller (2015) with 5,000 replications, which is recommended given our number of clusters (12) for regression. The Algo (Fired) Dummy takes value one for the *Algo* treatment (when a worker has been fired) and value zero otherwise. Ability is the performance on the task completed by all participants at the beginning of the experiment.

Dependent Variable	Change in production	
	Bias tasks [1]	Handicap tasks [2]
Constant	20.034*** (7.299)	51.581*** (14.177)
Algo Dummy	-5.263** (2.284)	-8.219 (11.590)
Fired Dummy	-29.059*** (6.760)	-66.798*** (16.588)
Algo × Fired Dummy	23.070*** (6.617)	46.673*** (17.554)
Period	-3.730*** (0.750)	-7.694*** (2.110)
Ability	-0.065 (0.084)	-0.070 (0.170)
Male Dummy	-9.626** (4.226)	-31.942*** (10.561)
N	316	316
R ²	0.153	0.117
Prob > χ^2	< 0.001	< 0.001
<u>P-values (F-Tests)</u>		
Algo × Fired + Algo	0.012	0.011
Algo × Fired + Fired	<0.001	0.005

*** Significant at the 0.01 level; ** at the 0.05 level; \diamond for p -values in (0.05, 0.10).

Attitudes

We report that for bias tasks ‘Algo × Fired’ is positive yet not significant (see regression [1] in Table B4, p -value = 0.097). The *Algo* treatment reduced perceived distributive justice for non-fired workers (see ‘Algo’, p -value = 0.029) and increased it for fired workers though not significantly so (see ‘Algo × Fired’ + ‘Algo’ = 0, p -value = 0.283). We note that the total effect of algorithm dismissals on perceived distributive justice is not significant (p -value = 0.943). Finally, perceived distributive justice was significantly lower for fired workers in the *Human* (see ‘Fired’, p -value < 0.001) and *Algo* treatments (see ‘Algo × Fired’ + ‘Fired’ = 0, p -value < 0.001).

Regarding work satisfaction, we report that ‘Algo × Fired’ is not significant (see regression [2] in Table B4, p -value = 0.840). The *Algo* treatment did not significantly impact work satisfaction for non-fired workers (see ‘Algo’, p -value = 0.840) and fired workers (see ‘Algo × Fired’ + ‘Algo’ = 0, p -value = 0.377). We note that the total effect of algorithm dismissals on work satisfaction is not significant (p -value = 0.314). Finally, work satisfaction was significantly lower for fired workers in the *Human* (see ‘Fired’, p -value = 0.009) and *Algo* treatments (see ‘Algo × Fired’ + ‘Fired’ = 0, p -value = 0.001).

4.3.3. Hypothesis 3

Performance

In handicap tasks, as expected, the worker inheriting the easier task produced more than the other worker in all cases. As a result, algorithms always dismissed the handicapped worker whereas human bosses did so in 71.8% of the cases.

In line with Hypothesis 3, the interaction coefficient ‘Algo × Fired’ is positive and significant (p -value = 0.008, see regression [2] in Table 3). Furthermore, the *Algo* treatment has a positive effect on the change in production for fired workers (‘Algo × Fired’ + ‘Algo’ = 0, p -value = 0.011) but not for non-fired workers (see ‘Algo’, p -value = 0.478). We note that the total effect of algorithm dismissals on the change in production is positive yet non-significant (p -value = 0.132). Finally, we find that for handicap tasks, the change in production is significantly lower for fired workers in the *Human* (see ‘Fired’, p -value < 0.001) and *Algo* treatments (see ‘Algo × Fired’ + ‘Fired’ = 0, p -value = 0.005).

Attitudes

In line with Hypothesis 3, we report that for *handicap* tasks ‘Algo × Fired’ is positive and significant (see regression [3] in Table B4 in Appendix B, p -value < 0.001). The *Algo* treatment reduced perceived distributive justice for non-fired workers (see ‘Algo’, p -value = 0.015) but not for fired workers (see ‘Algo × Fired’ + ‘Algo’ = 0, p -value = 0.661). We note that the total effect of algorithm dismissals on perceived distributive justice is negative yet not significant (p -value = 0.165). Finally, perceived distributive justice was significantly lower for fired workers in the *Human* (see ‘Fired’, p -value < 0.001) and *Algo* treatments (see ‘Algo × Fired’ + ‘Fired’ = 0, p -value < 0.001).

Regarding work satisfaction, we report that ‘Algo × Fired’ is also positive and significant (see regression [4] in Table B4, p -value = 0.004). The *Algo* treatment increased work satisfaction significantly for fired workers (see ‘Algo × Fired’ + ‘Algo’ = 0, p -value = 0.032) and decreased it for non-fired workers although not significantly so (see ‘Algo’, p -value = 0.256). We note that the total effect of algorithm dismissals on work satisfaction is not significant (p -value = 0.561). Finally, work satisfaction was significantly lower for fired workers in the *Human* treatment (see ‘Fired’, p -value < 0.001) but not in the *Algo* treatment (see ‘Algo × Fired’ + ‘Fired’ = 0, p -value = 0.213).

4.3.4. Hypothesis 4

By contrast with Hypothesis 4, we do not report a significant moderating role of algorithm aversion, as measured with the 5-item NARS (Syrdal et al., 2009). Indeed, the triple interaction (‘Algo × Fired × NARS’) and (‘Algo × Fired × NARS’ + ‘Algo × NARS’) are not significant, regardless of the dependent variable (see regressions [1], [2] and [3] in Table B5 in Appendix B). Thus, the NARS scores of workers did not impact the effect of algorithmic dismissals. These results could be due to the limited reliability of the scale (Cronbach α = 0.698) (Gliem and Gliem, 2003) or to the fact that our *Algo* treatment did not trigger substantial algorithm aversion. Algorithm aversion might have been reduced by our design choice that favored a simple algorithm which is both transparent and deprived of apparent biases. We could also speculate that people who have high NARS scores are also those who might set high standards regarding the fairness of human decisions. It might thus be that high-NARS workers are also easily disappointed by the decisions of human bosses.

4.4. Discussion

Our findings emphasize an asymmetry in the direction and magnitude of the reaction of fired and non-fired workers to algorithmic dismissals. The negative response of fired workers to the use of human versus algorithmic decisions is intense whereas the positive reaction of non-fired workers is mild. These results could be explained by an asymmetric reaction of workers to gains and losses (Kahneman and Tversky, 1979). Dismissed workers experience losses as their pay goes down, which tends to trigger a more intense response than that of workers who keep their job (e.g., Keysar et al., 2008). The negative reaction of fired workers could result from negative reciprocity (Fehr, Gächter and Kirchsteiger, 1997; Bolton and Ockenfels, 2000; Fehr and Gächter, 2000; Orhun, 2018). Moreover, negative reciprocity will tend to be magnified in the case of human decisions

because they are perceived, unlike algorithmic ones, as intentional (Falk, Fehr and Fischbacher, 2008). In the next section, we test the explanatory power of negative reciprocity concerns.

4.5. Spite hypothesis

As an exploratory, not pre-registered, hypothesis we test a key behavioral mechanism underlying our hypotheses. In the hypotheses section, we emphasized that the effort of fired workers in the second stage was partly driven by their resentment towards the human boss whose decisions were perceived as intentional. Following Falk, Fehr and Fischbacher (2008), we posit that perceived intentions behind human bosses' decisions matter when workers evaluate the fairness of dismissals. If unfair dismissal decisions are perceived to be intentional, workers' emotional reaction will be stronger (Lee, 2018) thus triggering spiteful behavior (Ben-Shakhar et al., 2007; Hopfensitz and Reuben, 2009; Bolle, Tan and Zizzo, 2014; Dickinson and Masclet, 2015). This will induce lower effort in the *Human* than in the *Algo* treatment in Stage 2. By contrast, workers who are not dismissed will perceive the tenure decision more positively in the *Human* than in the *Algo* treatment and will tend to exert more effort in the second stage. Our exploratory hypothesis can thus be stated as follows.

Hypothesis SP (Spite)

For all tasks, the positive *Algo treatment effect* will be more pronounced for fired workers who are spiteful. The positive difference in the *Algo treatment effect* between fired and non-fired workers will be more pronounced for those who are spiteful.

To test Hypothesis *SP*, we use the social preferences test developed in Bartling et al., (2009) and extended by Corgnet, Espín and Hernán-González (2015) (see A.2 in Appendix A). Spitefulness is defined as the number of decisions (out of 6) in which participants show that they value their relative standing. Spitefulness thus closely relates to envy (Bartling et al., 2009) that is reflected by an aversion to earning less than the other participant. For example, in Decision 3, spiteful individuals would choose the equal sharing of payoffs (10 euros each) rather than an allocation giving 18 euros to the other participant and only 10 euros to themselves. Following Corgnet, Espín and Hernán-González (2015), we define as spiteful a person who selected the spiteful option in

most of the decisions, that is in at least 4 out of 6 decisions. We thus categorize 28.5% of workers as spiteful.²⁰

In line with Hypothesis *SP*, we show in Table 4 (regression [1]), that the difference in the decrease in production of fired workers between treatments is significantly more pronounced for spiteful workers.²¹ This is the case because the triple interaction term ‘Algo × Fired × Spiteful’ is positive and significant (p -value = 0.042). The impact of spitefulness of fired workers on the effect of algorithmic dismissals is also positive, yet not significant (see ‘Algo × Fired × Spiteful’ + ‘Algo × Spiteful’, p -value = 0.065). The term ‘Fired × Spiteful’ is also negative and significant, thus showing that spiteful workers react more negatively to being fired than other workers. In line with Hypothesis *SP*, we find similar results for perceived distributive justice and work satisfaction (see regressions [2] and [3]) although the triple interaction ‘Algo × Fired × Spiteful’ is not significant for the latter.

Table 4. Change in production, perceived distributive justice, and work satisfaction as a function of treatments and spitefulness.

This table presents the results of panel regressions with random effects and clustered standard errors (in parentheses) at the session level using wild bootstrapping techniques as in Cameron and Miller (2015) with 5,000 replications, which is recommended given our number of clusters (12). The Algo (Fired) [Spiteful] Dummy takes value one for the *Algo* treatment (when a worker has been fired) [when a worker is categorized as spiteful] and value zero otherwise. Ability is the performance on the task completed by all participants at the beginning of the experiment. We include fixed effects for task types.

Dependent Variable	Change in production [1]	Perceived distributive justice [2]	Work satisfaction [3]
Constant	29.455*** (10.623)	16.090*** (0.529)	19.991*** (1.005)
Algo Dummy	-6.489 (6.309)	-0.989** (0.493)	-0.668 (0.819)
Fired Dummy	-34.030*** (6.124)	-4.862*** (0.380)	-2.670*** (0.421)
Algo × Fired	21.946*** (8.194)	1.186*** (0.446)	1.239** (0.618)
Algo × Fired × Spiteful	44.278** (21.797)	2.522** (1.278)	0.141 (1.483)
Algo × Spiteful	-2.608 (13.011)	-0.517 (0.979)	2.990 (2.072)

²⁰ This estimate is similar to the one reported in Capraro et al., (2017) (24.6%, Proportion Test, p -value = 0.349), using the same definition and a large dataset of US participants ($n = 370$).

²¹ We consider all seven periods as in Table 1.

Fired × Spiteful	-31.838** (16.241)	-1.841*** (0.666)	-0.286 (1.377)
Spiteful Dummy	-6.580* (3.583)	0.838 (0.770)	-1.033 (1.925)
Period	-5.679*** (1.366)	0.108** (0.055)	-0.242*** (0.071)
Ability	-0.050 (0.059)	-0.003 (0.005)	0.011 (0.008)
Male Dummy	-17.628*** (4.829)	-0.632 \diamond (0.332)	-1.796*** (0.614)
N	1,106	1,106	1,106
R ²	0.108	0.425	0.116
Prob > χ^2	<0.001	<0.001	<0.001
<u>P-values (F-Tests)</u>			
Algo × Fired × Spiteful + Algo × Spiteful	0.065	0.060	0.074

*** Significant at the 0.01 level; ** at the 0.05 level; \diamond for p -values in (0.05, 0.10).

5. Conclusion

The use of algorithms in management decisions is growing quickly (De Cremer, 2020) but our knowledge of the field continues to rely mostly on anecdotes and opinion pieces. Our goal was to start filling this gap in the literature by focusing on one of the most critical decisions of managers (Badaracco, 1997, 2016): dismissals. In the absence of archival data, we designed an experimental protocol to study workers' reactions to algorithmic demotion and provide a first set of guidelines to practitioners.

At a conceptual level, we developed and tested hypotheses regarding workers' behavioral and attitudinal reactions to algorithmic dismissals. We identified many situations in which algorithmic dismissals could alleviate the negative reaction of workers to firing decisions. This includes cases in which human dismissals are perceived as unfair because they are based on inaccurate performance metrics or biased by dimensions that are irrelevant to job performance. Our findings are clear-cut. Workers are less likely to react negatively to algorithmic than to human dismissals in a broad range of tasks. We show that this could be due to spiteful workers reacting less negatively to algorithmic than human dismissals. At a practical level, our findings imply that, unlike what is often heard in popular media, there are situations in which workers might embrace algorithmic

dismissals. At a legal level, our findings suggest GDPR regulators in Europe should contemplate settings in which fully automated decisions are allowed.

Because of the lack of a validated framework to study algorithmic dismissals, we developed a stylized experiment that captures essential features of firing decisions. Doing so, we left aside other relevant aspects of algorithms and the decision environment. An exciting avenue for future research on algorithmic dismissals will be to study automated rules that are not fully disclosed to workers. The decisions produced by these algorithms will likely be perceived by workers as less transparent and more biased (see Cowgill and Tucker, 2020). Using the current experimental design, we could deploy these opaque algorithms in a new set of tasks for which performance cannot be objectively measured (Castelo, Bos and Lehmann, 2019; Prah1 and Van Swol, 2021). The lack of transparency of algorithms and the use of tasks requiring subjective evaluations is likely to trigger more opposition to algorithmic dismissals (see Mahmud et al., 2022 for a review).

One should interpret our results as evidence that transparent algorithms can tame the negative reaction of fired workers to dismissals in a context in which there is a clear performance metric. This simple case is relevant because it captures the context of low-skill jobs that motivated our research (Lecher, 2019). Furthermore, the study of low-skill jobs is a pressing issue because these are the types of jobs companies will use to start experimenting with algorithmic dismissals.

The study of algorithmic dismissals necessarily calls for future research on algorithmic promotions. Our findings suggest that algorithms might be less effective at promoting than demoting workers. Because algorithms can tame the negative emotional reaction of workers to dismissals, they might also reduce the positive impact of promotion on a worker's morale. Our findings provide preliminary support for this claim because workers who kept their job perceived distributive justice to be lower when algorithms made the decision. Ultimately, the question that will haunt us is: can algorithms be effective leaders?

6. References

- Adams, J.S. (1965). "Inequality in social exchange". *Advanced Experimental Psychology*, 62: 335–343.
- Alekseev, A. (2020). *The Economics of Babysitting a Robot*. Available at SSRN 3656684.
- Badaracco, L. (1997). *Defining moments: When managers must choose between right and right*. Harvard Business Press.
- Badaracco, L. (2016). How to Tackle Your Toughest Decisions. *Harvard Business Review*, 104-107.
- Bai, B., Dai, H., Zhang, D., Zhang, F., & Hu, H. (In Press). The impacts of algorithmic work assignment on fairness perceptions and productivity: Evidence from field experiments. *Manufacturing & Service Operations Management*.
- Bartling, B., Fehr, E., Maréchal, M. A., & Schunk, D. (2009). Egalitarianism and competitiveness. *American Economic Review*, 99(2), 93-98.
- Ben-Shakhar, G., Bornstein, G., Hopfensitz, A., & Van Winden, F. (2007). Reciprocity and emotions in bargaining using physiological and self-report measures. *Journal of Economic Psychology*, 28(3), 314-323.
- Bigman, Y. E., & Gray, K. (2018). People are averse to machines making moral decisions. *Cognition*, 181, 21-34.
- Blau, F. D., & Kahn, L. M. (2017). The gender wage gap: Extent, trends, and explanations. *Journal of Economic Literature*, 55(3), 789-865.
- Blount, S. (1995). When social outcomes aren't fair: The effect of causal attributions on preferences. *Organizational Behavior and Human Decision Processes*, 63(2), 131-144.
- Bogert, E., Schechter, A., & Watson, R. T. (2021). Humans rely more on algorithms than social influence as a task becomes more difficult. *Scientific Reports*, 11(1), 1-9.

- Bolle, F., Tan, J. H., & Zizzo, D. J. (2014). Vendettas. *American Economic Journal: Microeconomics*, 6(2), 93-130.
- Bolton, G. & Ockenfels A. 2000. ERC: A theory of equity, reciprocity, and competition. *American Economic Review*, 90(1): 166-193.
- Bolton, G. E., & Ockenfels, A. (2008). Self-centered fairness in games with more than two players. *Handbook of Experimental Economics Results*, 1, 531-540.
- Bejarano, H., Corgnet, B., & Gómez-Miñambres, J. (2021). Economic stability promotes gift-exchange in the workplace. *Journal of Economic Behavior & Organization*, 187, 374-398.
- Brynjolfsson, E. & McAfee, A. (2014). *The second machine age: Work, progress, and prosperity in a time of brilliant technologies*. W. W. Norton & Company.
- Cameron, A. C., & Miller, D. L. (2015). A practitioner's guide to cluster-robust inference. *Journal of Human Resources*, 50(2), 317-372.
- Capraro, V., Corgnet, B., Espín, A. M., & Hernán-González, R. (2017). Deliberation favours social efficiency by making people disregard their relative shares: evidence from USA and India. *Royal Society Open Science*, 4(2), 160605.
- Charness, G. & Rabin, M. 2002. Understanding social preferences with simple tests. *Quarterly Journal of Economics*, 117: 817-869.
- Castelo, N., Bos, M. W., & Lehmann, D. R. (2019). Task-dependent algorithm aversion. *Journal of Marketing Research*, 56(5), 809-825.
- Chien, S. E., Chu, L., Lee, H. H., Yang, C. C., Lin, F. H., Yang, P. L., ... & Yeh, S. L. (2019). Age difference in perceived ease of use, curiosity, and implicit negative attitude toward robots. *ACM Transactions on Human-Robot Interaction (THRI)*, 8(2), 1-19.
- Chugunova, M., & Sele, D. (2020). We and it: An interdisciplinary review of the experimental evidence on human-machine interaction. *Max Planck Institute for Innovation & Competition Research Paper*, (20-15).

Colquitt, J. A. (2001). On the dimensionality of organizational justice: a construct validation of a measure. *Journal of Applied Psychology*, 86(3), 386.

Colquitt, J. A., Scott, B. A., Rodell, J. B., Long, D. M., Zapata, C. P., Conlon, D. E., & Wesson, M. J. (2013). Justice at the millennium, a decade later: a meta-analytic test of social exchange and affect-based perspectives. *Journal of Applied Psychology*, 98(2), 199.

Corgnet, B., Espín, A. M., & Hernán-González, R. (2015). The cognitive basis of social behavior: cognitive reflection overrides antisocial but not always prosocial motives. *Frontiers in Behavioral Neuroscience*, 9, 287.

Corgnet, B., Hernán-González, R., & Mateo, R. (2019). Race against the machine? Social incentives when humans meet robots. Working Paper. GATE 2019-04.

Cowgill, B. (2018). Bias and productivity in humans and algorithms: Theory and evidence from resume screening. Columbia Business School, Columbia University, 29.

Cowgill, B., & Tucker, C. E. (2020). Algorithmic fairness and economics. Columbia Business School Research Paper.

Crandall, J. W., Oudah, M., Ishowo-Oloko, F., Abdallah, S., Bonnefon, J. F., Cebrian, M., Shariff, A., Goodrich, M. A. & Rahwan, I. (2018). Cooperating with machines. *Nature Communications*, 9(1), 1-12.

Cross, E. S., & Ramsey, R. (2021). Mind meets machine: towards a cognitive science of human-machine interactions. *Trends in Cognitive Sciences*, 25(3), 200-212.

Daugherty, P. R., & Wilson, H. J. (2018). *Human + machine: Reimagining work in the age of AI*. Harvard Business Press.

De Cremer, D. (2020). *Leadership by algorithm: Who leads and who follows in the AI era?* Harriman House Limited.

Dickinson, D. L., & Masclet, D. (2015). Emotion venting and punishment in public good experiments. *Journal of Public Economics*, 122, 55-67.

- Dietvorst, B. J., Simmons, J. P., & Massey, C. (2015). Algorithm aversion: People erroneously avoid algorithms after seeing them err. *Journal of Experimental Psychology: General*, 144(1), 114.
- Dietvorst, B. J., Simmons, J. P., & Massey, C. (2018). Overcoming algorithm aversion: People will use imperfect algorithms if they can (even slightly) modify them. *Management Science*, 64(3), 1155-1170.
- Falk, A., Fehr, E., & Fischbacher, U. (2008). Testing theories of fairness—Intentions matter. *Games and Economic Behavior*, 62(1), 287-303.
- Fehr, E. & Fischbacher, U. 2002. Why social preferences matter – The impact of non-selfish motives on competition, cooperation, and incentives. *Economic Journal*, 112: C1-C33.
- Fehr, E., Gächter, S., & Kirchsteiger, G. (1997). Reciprocity as a contract enforcement device: Experimental evidence. *Econometrica*, 833-860.
- Fehr, E., & Gächter, S. (2000). Cooperation and punishment in public goods experiments. *American Economic Review*, 90(4), 980-994.
- Fehr, E. & Schmidt, K. 1999. A theory of fairness, competition, and cooperation. *The Quarterly Journal of Economics*, 114(3): 817-868.
- Filiz, I., Judek, J. R., Lorenz, M., & Spiwoкс, M. (2021). The tragedy of algorithm aversion. Ostfalia Hochschule für Angewandte Wissenschaften, Fakultät Wirtschaft.
- Fortson, D. 2021. Fired by AI: The algorithm that judges whether staff are really working from home. *The Times*, December 19th 2021.
<https://www.thetimes.co.uk/article/computer-says-clear-your-desk-artificial-intelligence-wfh-covid-hc7wvs3s3>
- Frey, C. B., & Osborne, M. A. (2017). The future of employment: How susceptible are jobs to computerization? *Technological forecasting and social change*, 114, 254-280.
- Fumagalli, E., Rezaei, S., & Salomons, A. (2022). OK computer: Worker perceptions of algorithmic recruitment. *Research Policy*, 51(2), 104420.

Gächter, S., & Riedl, A. (2005). Moral property rights in bargaining with infeasible claims. *Management Science*, 51(2), 249-263.

Gazzaniga, M. S. (2000). *Human-The Science Behind What Makes Us Unique*. Harper Perennial.

Gliem, J. A., & Gliem, R. R. (2003). Calculating, interpreting, and reporting Cronbach's alpha reliability coefficient for Likert-type scales. *Midwest Research-to-Practice Conference in Adult, Continuing, and Community Education*.

Gogoll, J., & Uhl, M. (2018). Rage against the machine: Automation in the moral domain. *Journal of Behavioral and Experimental Economics*, 74, 97-103.

Goldman, B. M. (2003). The application of referent cognitions theory to legal-claiming by terminated workers: The role of organizational justice and anger. *Journal of Management*, 29(5), 705-728.

Granulo, A., Fuchs, C., & Puntoni, S. (2019). Psychological reactions to human versus robotic job replacement. *Nature Human Behaviour*, 3(10), 1062-1069.

Haidt, J. (2012). *The righteous mind: Why good people are divided by politics and religion*. Vintage.

Hawking, S., Russell, S., Tegmark, M., & Wilczek, F. (2014). Transcendence looks at the implications of artificial intelligence – but are we taking AI seriously enough?. *Independent*, May 1, 2014.

Hidalgo, C. A., Orghian, D., Canals, J. A., De Almeida, F., & Martín, N. (2021). *How humans judge machines*. MIT Press.

Highhouse, S. (2008). Stubborn Reliance on Intuition and Subjectivity in Employee Selection. *Industrial and Organizational Psychology*, 1, 333-342.

Hoffman, M., Kahn, L. B., & Li, D. (2018). Discretion in hiring. *The Quarterly Journal of Economics*, 133(2), 765-800.

Hopfensitz, A., & Reuben, E. (2009). The importance of emotions for the effectiveness of social punishment. *The Economic Journal*, 119(540), 1534-1559.

Jamison, J., Karlan, D., & Schechter, L. (2008). To deceive or not to deceive: The effect of deception on behavior in future laboratory experiments. *Journal of Economic Behavior & Organization*, 68(3-4), 477-488.

Kagel, J. H., Kim, C., & Moser, D. (1996). Fairness in ultimatum games with asymmetric information and asymmetric payoffs. *Games and Economic Behavior*, 13(1), 100-110.

Kahneman, D., & Tversky, A. (1979). Prospect Theory: An Analysis of Decision under Risk. *Econometrica*, 47(2), 263-292.

Keynes, J. M. (1930). Economic possibilities for our grandchildren. In *Essays in Persuasion* (New York: Harcourt Brace, 1932), 358-373.

Keysar, B., Converse, B. A., Wang, J., & Epley, N. (2008). Reciprocity is not give and take: Asymmetric reciprocity to positive and negative acts. *Psychological Science*, 19(12), 1280-1286.

Kirchkamp, O., & Strobel, C. (2019). Sharing responsibility with a machine. *Journal of Behavioral and Experimental Economics*, 80, 25-33.

Kleinberg, J., Ludwig, J., Mullainathan, S., & Sunstein, C. R. (2018). Discrimination in the Age of Algorithms. *Journal of Legal Analysis*, 10, 113-174.

Kleinberg, J., Ludwig, J., Mullainathan, S., & Sunstein, C. R. (2020). Algorithms as discrimination detectors. *Proceedings of the National Academy of Sciences*, 117(48), 30096-30100.

Knez, M. J., & Camerer, C. F. (1995). Outside options and social comparison in three-player ultimatum game experiments. *Games and Economic Behavior*, 10(1), 65-94.

Konow, J. (2003). Which is the fairest one of all? A positive analysis of justice theories. *Journal of Economic Literature*, 41(4), 1188-1239.

Konow, J. (2000). Fair shares: Accountability and cognitive dissonance in allocation decisions. *American Economic Review*, 90(4), 1072-1091.

Lambrecht, A., & Tucker, C. (2019). Algorithmic bias? An empirical study of apparent gender-based discrimination in the display of STEM career ads. *Management Science*, 65(7), 2966-2981.

Lecher, C. (2019). How Amazon automatically tracks and fires warehouse workers for 'productivity'. *The Verge*, 25, 2019.

Lee, M. K. (2018). Understanding perception of algorithmic decisions: Fairness, trust, and emotion in response to algorithmic management. *Big Data & Society*, 5(1), 2053951718756684.

Leventhal, G. S., & Michaels, J. W. (1971). Locus of cause and equity motivation as determinants of reward allocation. *Journal of Personality and Social Psychology*, 17(3), 229.

Lind, E. A., Greenberg, J., Scott, K. S., & Welchans, T. D. (2000). The winding road from employee to complainant: Situational and psychological determinants of wrongful-termination claims. *Administrative Science Quarterly*, 45(3), 557-590.

Mahmud, H., Islam, A. N., Ahmed, S. I., & Smolander, K. (2022). What influences algorithmic decision-making? A systematic literature review on algorithm aversion. *Technological Forecasting and Social Change*, 175, 121390.

Manyika, J., & Sneider, K. (2018). AI, automation, and the future of work: Ten things to solve for. McKinsey Global Institute: June. Available at <https://www.mckinsey.com/featured-insights/future-of-work/ai-automation-and-the-future-of-work-ten-things-to-solve-for>.

March, C. (2021). Strategic interactions between humans and artificial intelligence: Lessons from experiments with computer players. *Journal of Economic Psychology*, 87, 102426.

McKee-Ryan, F., Song, Z., Wanberg, C. R., & Kinicki, A. J. (2005). Psychological and physical well-being during unemployment: a meta-analytic study. *Journal of Applied Psychology*, 90(1), 53.

Newman, D. T., Fast, N. J., & Harmon, D. J. (2020). When eliminating bias isn't fair: Algorithmic reductionism and procedural justice in human resource decisions. *Organizational Behavior and Human Decision Processes*, 160, 149-167.

Nozick, Robert. 1974. *Anarchy, State, and Utopia*. NY: Basic Books.

Orhun, A. Y. (2018). Perceived motives and reciprocity. *Games and Economic Behavior*, 109, 436-451.

Paul, K. I., & Moser, K. (2009). Unemployment impairs mental health: Meta-analyses. *Journal of Vocational Behavior*, 74(3), 264-282.

Prahl, A., & Van Swol, L. M. (2021). Out with the humans, in with the machines?: Investigating the behavioral and psychological effects of replacing human advisors with a machine. *Human-Machine Communication*, 2, 209-234.

Rauchbauer, B., Nazarian, B., Bourhis, M., Ochs, M., Prévot, L., & Chaminade, T. (2019). Brain activity during reciprocal social interaction investigated using conversational robots as control condition. *Philosophical Transactions of the Royal Society B*, 374(1771), 20180033.

Raveendhran, R., & Fast, N. J. (2021). Humans judge, algorithms nudge: The psychology of behavior tracking acceptance. *Organizational Behavior and Human Decision Processes*, 164, 11-26.

Riek, L. D. (2012). Wizard of Oz studies in HRI: A systematic review and new reporting guidelines. *Journal of Human-Robot Interaction*, 1(1), 119-136.

Ruffle, B. J. (1998). More is better, but fair is fair: Tipping in dictator and ultimatum games. *Games and Economic Behavior*, 23(2), 247-265.

Russell, S. (2019). *Human compatible: Artificial intelligence and the problem of control*. Penguin.

Russell, S., & Norvig, P. (2021). *Artificial intelligence: a modern approach, global edition 4th*. Foundations, 19, 23.

Ryan, R. M. (1982). Control and information in the intrapersonal sphere: An extension of cognitive evaluation theory. *Journal of Personality and Social Psychology*, 43(3), 450-461.

Schniter, E., Shields, T. W., & Sznycer, D. (2020). Trust in humans and robots: Economically similar but emotionally different. *Journal of Economic Psychology*, 78, 102253.

Shin, D., & Park, Y. J. (2019). Role of fairness, accountability, and transparency in algorithmic affordance. *Computers in Human Behavior*, 98, 277-284.

Simon, H. A. (1965). *The shape of automation for men and management* (Vol. 13). New York: Harper & Row.

Smith, T. (2020). Fired by robots — Uber faces legal challenge for algorithmic dismissals. Sifted. <https://sifted.eu/articles/uber-algorithm-firing-drivers/>

Stanovich, K. E. (2005). *The robot's rebellion: Finding meaning in the age of Darwin*. University of Chicago press.

Strobel, C. (2021). *The hidden costs of automation*. Working Paper. Hamburg University of Technology.

Syrdal, D. S., Dautenhahn, K., Koay, K. L., & Walters, M. L. (2009). The negative attitudes towards robots scale and reactions to robot behaviour in a live human-robot interaction study. *Adaptive and Emergent Behaviour and Complex Systems*.

Theodossiou, I. (1998). The effects of low-pay and unemployment on psychological well-being: a logistic regression approach. *Journal of Health Economics*, 17(1), 85-104.

Traeger, M. L., Strohkorb Sebo, S., Jung, M., Scassellati, B., & Christakis, N. A. (2020). Vulnerable robots positively shape human conversational dynamics in a human-robot team. *Proceedings of the National Academy of Sciences*, 117(12), 6370-6375.

Venn, D. (2009). *Legislation, collective bargaining and enforcement: updating the OECD employment protection indicators*. OECD social, employment and migration working papers 89.

Wang, L., Rau, P. L. P., Evers, V., Robinson, B. K., & Hinds, P. (2010, March). When in Rome: the role of culture & context in adherence to robot recommendations. In *2010 5th ACM/IEEE International Conference on Human-Robot Interaction (HRI)* (pp. 359-366). IEEE.

Weichselbaumer, D., & Winter-Ebmer, R. (2005). A meta-analysis of the international gender wage gap. *Journal of Economic Surveys*, 19(3), 479-511.

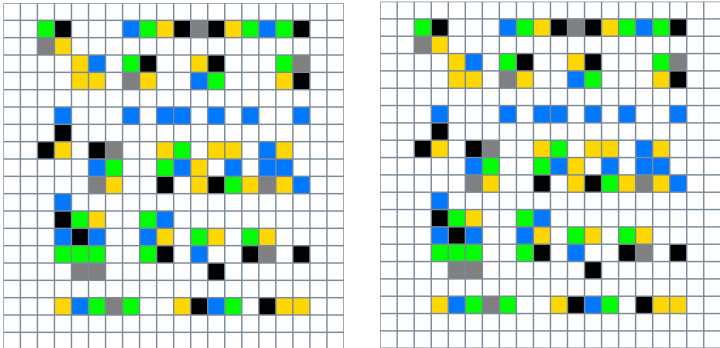
Appendices

Appendix A (Instructions)

A.1. Tasks

Worker 1

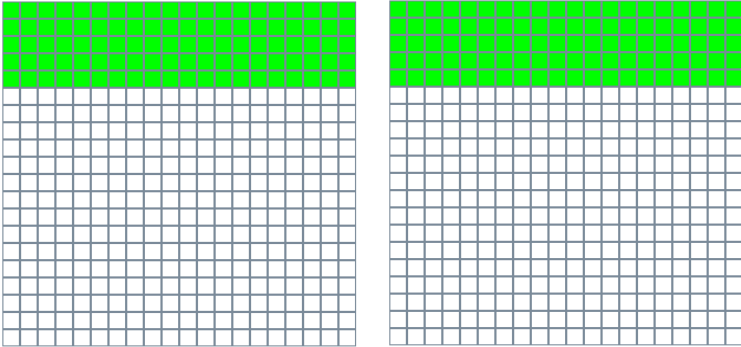
Worker 2



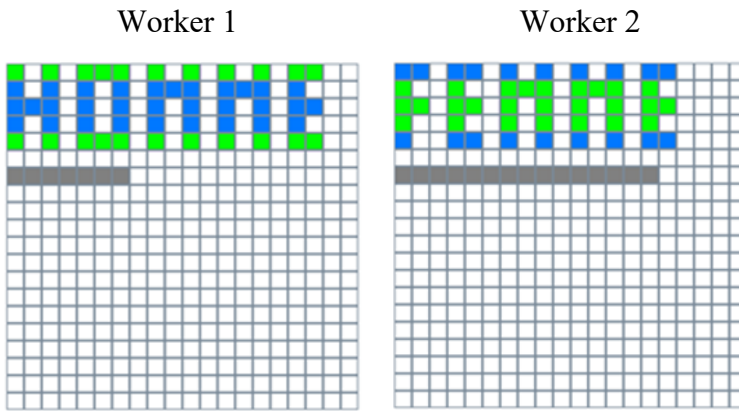
Ability task

Worker 1

Worker 2



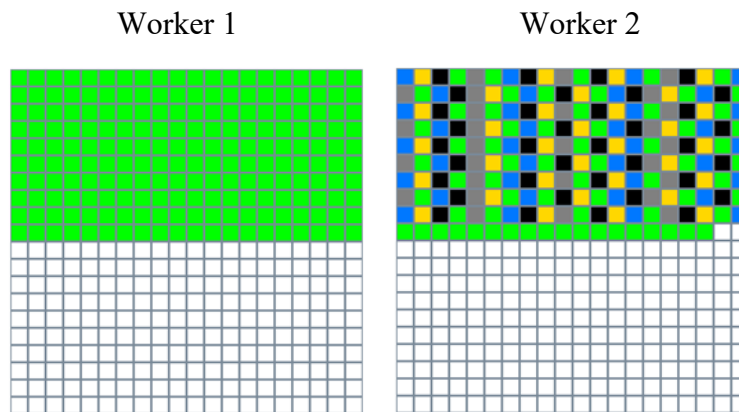
Task 1. Easy measurable task



Task 2. *Bias* task I

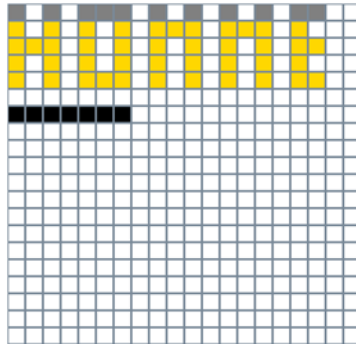


Task 3. *Hard measurable* task

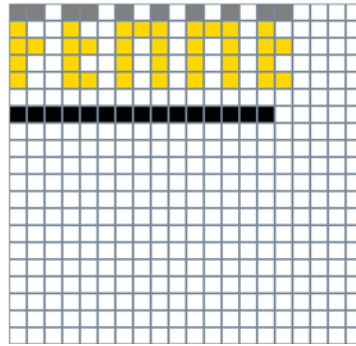


Task 4. *Handicap* task I

Worker 1



Worker 2

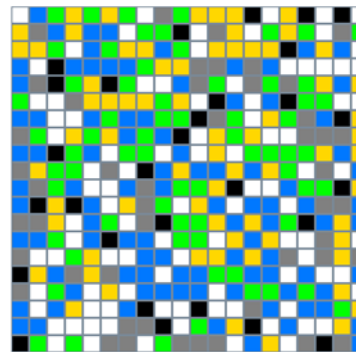


Task 5. *Bias* task II

Worker 1

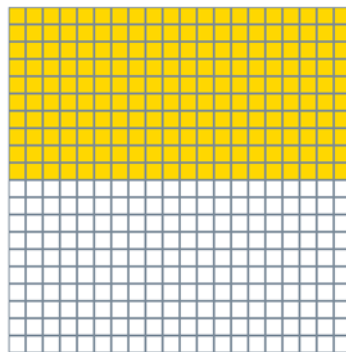


Worker 2

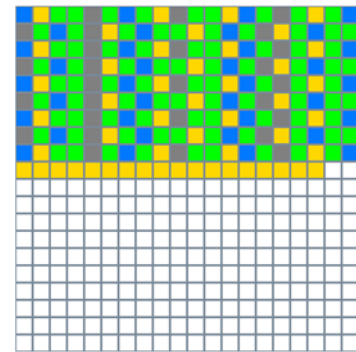


Task 6. *Non-measurable* task

Worker 1



Worker 2



Task 7. *Handicap* task II

A.2. Translated instructions

We highlight in *italics* the parts of the instructions that apply to the *Algo* treatment.

Welcome screen

Thank you for participating in this experiment on decision making. Please turn off your cell phone and do not communicate with other participants until the end of the session.

If you have any questions, you can raise your hand or press the red button on the side of your desk at any time. We will come to answer you individually. During this session you will make several decisions. These decisions can earn you money. Regardless of these decisions, you will receive 5 euros for showing up on time.

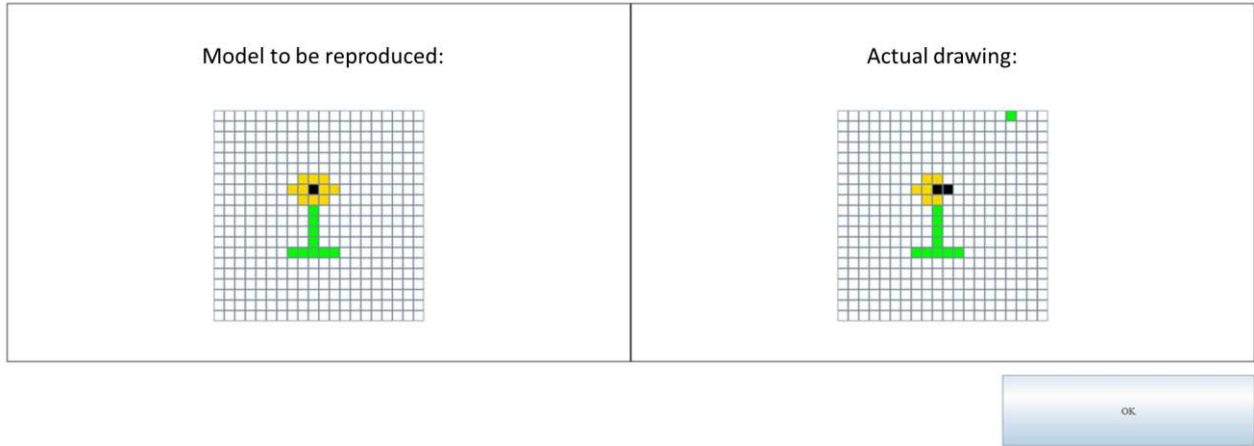
This experiment is composed of seven periods of the main task and a final questionnaire.

All instructions are on your screen. Please click OK to continue

Task presentation

This task consists of reproducing a model as accurately as possible on a grid composed of 400 cells. The pay obtained for this work will depend directly on your performance. The value of each of the cells in the grid is equal to 2 cents if it is filled with the same color as the model. The value of the cell is -1 cent if it is filled with a color other than the model. The value of the cell is 0 if it is not filled with any color.

In the following example, the model consists of 19 cells. The value of the completed grid is equal to 30 cents since 16 cells were filled with the same color as the model while 2 cells were filled with a wrong color. This is the case of the yellow cell in the center of the grid and the green cell at the top right of the grid. Note, however, that the minimum value of your work is set at 0 and therefore cannot be negative. You will have 90 seconds to reproduce the model.



Ability task screen

Time remaining

Please, reproduce the following model.
To draw, select a color box on the top of the grid and click on the cells you want to color.

Main phase of the experiment

Screen I

Each of the seven periods of the experiment consists of two stages.

Each participant will be randomly assigned to a role: P1, P2 or B. You will keep your role for the seven periods. If you are person B you will receive an endowment of 16 euros. P1 and P2 will not

receive any initial endowment. The initial endowment of B will be used to pay P1 and P2 during the two periods of a given phase.

At the beginning of each period, you will be matched with two other participants who are selected at random in this experimental session.

Screen II

During the first period, P1 and P2 will work on a task similar to the one that all the participants have just done. Participant B will have to pay 5 euros each of the workers (P1 and P2) regardless of their level of performance on the task. P1 and P2 will only get this fixed payment. They will not be paid based on their performance. On the other hand, participant B will be paid according to the exact value of the work done by the two workers P1 and P2.

At the end of the first period, B *<a pre-programmed robot>* will decide which of the two workers P1 or P2 will be kept at his workstation for the second period. The worker kept in his or her position will continue to be paid a fixed payment of 5 euros to carry out the task while the other person will only receive 1 euro of fixed payment. Both workers will however be able to perform the task and therefore generate value for participant B.

Screen III

To make his or her decision, Participant B will see a screen representing each worker's grid and model as well as the exact value generated by each worker. This information screen will be displayed for exactly 25 seconds after which the participant will have to make their decision. Workers will receive the same information on their screen and will be informed of which worker is maintained and which is demoted.

<To make its decision the robot will compute the value of the grid produced by each worker. The worker who produces the most value will be maintained, and the other worker will be demoted. In case of a tie, the robot selects a worker at random. Participant B and the two workers will see a screen representing each worker's grid and model as well as the exact value generated by each worker. This information screen will be displayed for exactly 25 seconds.>

Earnings screen

For the calculation of your final earnings on the main phase of the experiment, a period (out of 7) will be chosen at random.

The earnings for a given period are calculated as follows depending on your role:

You have been assigned the role of P1 or P2:

And you were maintained at your workstation at the end of the first period:

Earnings = 10 euros (5 euros for the first period + 5 euros for the second period)

And you were demoted at the end of the first period:

Earnings = 6 euros (5 euros for the first period + 1 euro for the second period)

You have been assigned role B:

16 euros of initial endowment – 16 euros (payments to participants P1 and P2 in both periods) + value of the two grids completed by P1 and P2 in both periods.

To these earnings will be added your earnings for the initial task, the final questionnaire and the comprehension quiz.

Comprehension quiz

Please, answer the follow 7 questions correctly to continue the experiment. At the end of the experiment, you will receive a 1-euro bonus if you answered all questions correctly in your first attempt. <in **bold** is the correct answer> [in *italics* for the *Algo* treatment]

1. What is the structure of experience?

a) There are two six-minute periods; b) There are six one-minute periods; c) There are seven periods divided into two stages; d) There are two periods of ten minutes; e) There is one 20- minute period.

2. How are roles determined in the experiment?

a) The role of a participant (P1, P2 or B) is randomly determined at the start of each period; b) The role of a participant (P1, P2 or B) is randomly determined at the start of each stage; c) Participants choose their role (P1, P2 or B); d) The role of a participant (P1, P2 or B) is randomly determined at the start of the first period and remains the same for the entire experiment; e) The roles (P1, P2 or B) are chosen according to the profile of each participant.

3. What is the value (in euro cents) of a drawing produced by a worker?

a) It is 3 times the number of cells colored correctly; b) It is 2 times the number of cells colored correctly minus the number of cells colored incorrectly; c) It is the number of cells colored

correctly; d) It is 4 times the number of cells colored correctly minus 2 times the number of cells colored incorrectly; e) It is the number of cells colored correctly minus the number of cells colored incorrectly.

4. Who decides which of the two workers (P1 or P2) is demoted in the second period?

a) A pre-programmed robot that decides randomly; b) Participant B chooses whether himself or a pre-programmed robot will make this decision; c) The workers themselves; d) A pre-programmed robot that demotes the worker who has produced the least; e) Participant B after observing the workers' drawings for 25 seconds.

5. How many periods of the experiments will be paid?

a) The first period; b) One of the seven periods chosen at random; c) The last period; d) All periods; e) Periods 2 and 6.

6. How are workers P1 and P2 compensated for the period chosen at random for the payments?

a) They get 4 euros for each of the two periods; b) They get 5 euros for the first period and 1 euro for the second period; c) They get 5 euros for both periods if they are maintained, and 5 euros for the first period and 1 euro for the second if they are demoted; d) They get 5 euros for the two periods if they are maintained, and 1 euro for the two periods if they are demoted; e) They get 1 euro for the first period and 5 euros for the second period.

7. How is Participant B compensated for the period chosen at random for the payments?

a) He or she gets 5 euros for each of the two periods.

b) He or she obtains the value of the drawing produced by the two workers in the two periods.

c) He or she only obtains the value of the drawing produced by the worker who has been maintained.

d) He or she obtains the value of the drawing produced by the two workers in the first period.

e) He or she obtains the value of the drawing produced by the two workers for one of the two periods chosen at random.

Questions at the end of each period

Perceived distributive justice (adapted from Colquitt, 2001)

Please read each statement and select the appropriate number, based on the following scale: 5 = very much, 4 = somewhat, 3 = neutral, 2 = not much, 1 = not at all.

Does your pay reflect the effort you put into your work?

Does your pay match the work you have done?

Is your pay justified, given your performance?

Work satisfaction (adapter from Ryan, 1982)

Please read each statement and select the appropriate number, based on the following scale: 7 = strongly agree, 6 = agree, 5 = somewhat agree, 4 = neutral, 3 = somewhat disagree, 2 = disagree, 1 = strongly disagree.

I really enjoyed doing this task.

This task was fun to do.

I thought it was a boring task.

I found this task quite enjoyable.

Final Questionnaire

Negative Attitudes towards Robots Scale (NARS) (Syrdal et al., 2009)

Please read each statement and decide how much you agree or disagree with its content. Then select the appropriate number, based on the following scale: 5 = completely agree, 4 = agree, 3 = neutral (neither agree nor disagree), 2 = disagree, 1 = strongly disagree.

1- I would feel uncomfortable in a job where I had to use robots.

2- The word “robot” means nothing to me.

3- I hate the idea that robots or artificial intelligences can make judgments on people.

4- I would feel very nervous in front of a robot.

5- I would feel anxious if I had to talk to a robot.

Social preferences test (adapted from Bartling et al., 2009; Corgnet, Espín and Hernán-González, 2015)

In this part of the experiment, you will be asked to make a series of choices in decision problems. For each line in the table in the next screen, please state whether you prefer option A or option B. Notice that there are a total of 6 lines in the table but just one line will be randomly selected for payment. Each line is equally likely to be chosen, so you should pay equal attention to the choice you make in every line.

Your earnings for the selected line depend on which option you chose: if you chose option A in that line, you will receive 10 euros and the other participant who will be matched with you will also receive 10 euros. If you chose option B in that line, you and the other participant will receive earnings as indicated in the table for that specific line.

1	Option A: <input type="radio"/> 10€ for you <input type="radio"/> 10€ for the other participant	Option B: <input type="radio"/> 10€ for you <input type="radio"/> 6€ for the other participant
2	Option A: <input type="radio"/> 10€ for you <input type="radio"/> 10€ for the other participant	Option B: <input type="radio"/> 16€ for you <input type="radio"/> 4€ for the other participant
3	Option A: <input type="radio"/> 10€ for you <input type="radio"/> 10€ for the other participant	Option B: <input type="radio"/> 10€ for you <input type="radio"/> 18€ for the other participant
4	Option A: <input type="radio"/> 10€ for you <input type="radio"/> 10€ for the other participant	Option B: <input type="radio"/> 11€ for you <input type="radio"/> 19€ for the other participant
5	Option A: <input type="radio"/> 10€ for you <input type="radio"/> 10€ for the other participant	Option B: <input type="radio"/> 12€ for you <input type="radio"/> 4€ for the other participant
6	Option A: <input type="radio"/> 10€ for you <input type="radio"/> 10€ for the other participant	Option B: <input type="radio"/> 8€ for you <input type="radio"/> 16€ for the other participant

For example, if you chose B in line 2 and this line is selected for payment, you will receive euros 16 and the other participant will receive 4 euros. Similarly, if you chose B in line 3 and this line is selected for payment, you will receive 10 euros and the other participant will receive euros 18. Note that the other participant will never be informed of your personal identity, and you will not be informed of the other participant's personal identity.

After all of you have made their choices, the computer will select two and only two participants in the room. The decision table of the first participant will determine the payoff of the two subjects. Then the computer will randomly determine which line of the first subject decision table is going to be paid.

Demographics (age, gender and school)

Appendix B (Additional results)

Table B1. Production and stages.

This table presents the results of linear panel regressions with random effects and clustered standard errors (in parentheses) at the session level using wild bootstrapping techniques as in Cameron and Miller (2015) with 5,000 replications, which is recommended given our number of clusters (12). The *Algo* Dummy takes value one for the *Algo* treatment and value zero otherwise. Ability is the performance on the task completed by all participants at the beginning of the experiment. We include fixed effects for task types.

Dependent Variable	Production in Stage 1 [1]	Production in Stage 2 [2]
Constant	443.846*** (9.639)	447.561*** (16.195)
Algo Dummy	0.667 (4.178)	10.747 \diamond (5.932)
Period	3.370*** (0.334)	-2.342 \diamond (1.292)
Ability	0.174 \diamond (0.093)	0.171 (0.127)
Male Dummy	-4.968 (7.250)	-22.418** (10.593)
Number of Observations	1,106	1,106
R ²	0.746	0.585
Prob > χ^2	<0.001	<0.001

*** Significant at the 0.01 level; ** at the 0.05 level; \diamond for p -values in (0.05, 0.10).

Table B2. Perceived distributive justice and task measurability.

This table presents the results of panel <OLS> regressions with random effects and clustered standard errors (in parentheses) at the session level using wild bootstrapping techniques as in Cameron and Miller (2015) with 5,000 replications, which is recommended given our number of clusters (12) for regression [1] <[2], [3] and [4]>. The Algo (Fired) Dummy takes value one for the *Algo* treatment (when a worker has been fired) and value zero otherwise. Ability is the performance on the task completed by all participants at the beginning of the experiment. We include fixed effects for task type in regression [1].

Dependent Variable	Perceived distributive justice			
	Measurable tasks [1]	Easy & Measurable task [2]	Hard & Measurable task [3]	Non-measurable task [4]
Constant	16.105*** (0.476)	14.565*** (0.555)	15.600*** (0.623)	15.884*** (0.637)
Algo Dummy	-0.961** (0.455)	-0.913 \diamond (0.525)	-1.045 \diamond (0.623)	-0.704 (0.495)
Fired Dummy	-5.226*** (0.370)	-6.305*** (0.419)	-4.466*** (0.568)	-3.820*** (0.637)
Algo \times Fired Dummy	1.935*** (0.539)	1.459** (0.733)	2.462*** (0.738)	1.176 (0.801)
Ability	0.005 (0.005)	0.004 (0.006)	0.005 (0.007)	-0.004 (0.007)
Male Dummy	-0.749 \diamond (0.409)	-0.915** (0.416)	-0.754 (0.479)	-0.647 (0.545)
N	316	158	158	158
R ²	0.423	0.521	0.250	0.230
Prob > χ^2	< 0.001	< 0.001	< 0.001	< 0.001
<u>P-values (F-Tests)</u>				
Algo \times Fired + Algo	0.092	0.370	0.066	0.568
Algo \times Fired + Fired	< 0.001	< 0.001	< 0.001	< 0.001

*** Significant at the 0.01 level; ** at the 0.05 level; \diamond for *p*-values in (0.05, 0.10).

Table B3. Work satisfaction and task measurability.

This table presents the results of panel <OLS> regressions with random effects and clustered standard errors (in parentheses) at the session level using wild bootstrapping techniques as in Cameron and Miller (2015) with 5,000 replications, which is recommended given our number of clusters (12) for regression [1] <[2], [3] and [4]>. The Algo (Fired) Dummy takes value one for the *Algo* treatment (when a worker has been fired) and value zero otherwise. Ability is the performance on the task completed by all participants at the beginning of the experiment. We include fixed effects for task type in regression [1].

Dependent Variable	Work satisfaction			
	Measurable tasks [1]	Easy & Measurable task [2]	Hard & Measurable task [3]	Non-measurable task [4]
Constant	18.163*** (0.769)	20.650*** (1.158)	18.293*** (0.857)	18.801*** (1.130)
Algo Dummy	1.178 \diamond (0.612)	1.246 \diamond (0.698)	0.461 (1.011)	-0.603 (0.783)
Fired Dummy	-1.515** (0.720)	-2.961*** (0.806)	-2.096 (2.025)	-3.621*** (1.015)
Algo \times Fired Dummy	0.326 (0.962)	-0.048 (1.207)	1.969 (2.271)	3.201** (1.533)
Ability	0.012 (0.010)	0.011 (0.014)	0.009 (0.008)	0.001 (0.012)
Male Dummy	-2.091*** (0.722)	-3.046*** (0.727)	-1.328 (0.897)	-0.787 (0.823)
N	316	158	158	158
R ²	0.106	0.208	0.049	0.062
Prob > χ^2	< 0.001	< 0.001	< 0.001	< 0.001
P-values (F-Tests)				
Algo \times Fired + Algo	0.128	0.284	0.139	0.020
Algo \times Fired + Fired	0.071	0.001	0.915	0.731

*** Significant at the 0.01 level; ** at the 0.05 level; \diamond for *p*-values in (0.05, 0.10).

Table B4. Perceived distributive justice and work satisfaction for bias and handicap tasks.

This table presents the results of panel regressions with random effects and clustered standard errors (in parentheses) at the session level using wild bootstrapping techniques as in Cameron and Miller (2015) with 5,000 replications, which is recommended given our number of clusters (12). The Algo (Fired) Dummy takes value one for the *Algo* treatment (when a worker has been fired) and value zero otherwise. Ability is the performance on the task completed by all participants at the beginning of the experiment.

Dependent Variable	Perceived distributive justice	Work satisfaction	Perceived distributive justice	Work satisfaction
	Bias tasks [1]	Bias tasks [2]	Handicap tasks [3]	Handicap tasks [4]
Constant	14.458*** (0.706)	21.307*** (0.875)	14.129*** (0.607)	19.028*** (1.634)
Algo Dummy	-0.898** (0.411)	0.491 (0.869)	-1.518** (0.624)	-1.118 (0.985)
Fired Dummy	-5.782*** (0.529)	-2.598*** (0.996)	-5.356*** (0.613)	-3.923*** (0.926)
Algo × Fired Dummy	1.866 \diamond (1.126)	0.251 (1.243)	1.752** (0.731)	3.147*** (1.098)
Period	0.122 (0.087)	-0.346*** (0.113)	0.084 \diamond (0.045)	-0.137 (0.122)
Ability	-0.005 (0.006)	0.002 (0.010)	-0.006 (0.008)	0.026 \diamond (0.014)
Male Dummy	-0.329 (0.509)	-1.736*** (0.637)	-0.740 (0.482)	-1.950*** (0.592)
N	316	316	316	316
R ²	0.404	0.094	0.414	0.105
Prob > χ^2	< 0.001	< 0.001	< 0.001	< 0.001
<u>P-values (F-Tests)</u>				
Algo × Fired + Algo	0.283	0.377	0.661	0.032
Algo × Fired + Fired	< 0.001	0.001	< 0.001	0.213

*** Significant at the 0.01 level; ** at the 0.05 level; \diamond for p -values in (0.05, 0.10).

Table B5. Change in production, perceived distributive justice and work satisfaction as a function of treatments and NARS scores.

This table presents the results of panel regressions with random effects and clustered standard errors (in parentheses) at the session level using wild bootstrapping techniques as in Cameron and Miller (2015) with 5,000 replications, which is recommended given our number of clusters (12). The Algo (Fired) Dummy takes value one for the *Algo* treatment (when a worker has been fired) and value zero otherwise. Ability is the performance on the task completed by all participants at the beginning of the experiment. We include fixed effects for task types.

Dependent Variable	Change in production [1]	Perceived distributive justice [2]	Work satisfaction [3]
Constant	31.185 (26.155)	17.328*** (1.404)	20.309*** (4.751)
Algo × Fired × NARS	1.947 (3.124)	-0.067 (0.188)	-0.087 (0.205)
Algo × NARS	-0.470 (1.462)	0.235 [◊] (0.129)	0.082 (0.249)
Fired × NARS	-1.785 (1.462)	-0.011 (0.166)	0.030 (0.185)
NARS	-0.194 (1.048)	-0.057 (0.075)	-0.034 (0.238)
Algo Dummy	2.171 (26.500)	-5.450** (2.312)	-1.307 (4.528)
Fired Dummy	-10.866 (51.253)	-5.185 [◊] (2.816)	-3.276 (3.490)
Algo × Fired Dummy	-0.734 (56.951)	3.130 (3.166)	2.882 (3.880)
Period	-5.850 (1.367)	0.103 [◊] (0.059)	-0.241*** (0.069)
Ability	-0.046 (0.063)	-0.002 (0.004)	0.011 (0.008)
Male Dummy	-17.480*** (5.680)	-0.592 (0.361)	-1.818** (0.752)
N	1,106	1,106	1,106
R ²	0.094	0.421	0.100
Prob > χ^2	<0.001	<0.001	<0.001
<u>P-values (F-Tests)</u>			
Algo × Fired × NARS + Algo × NARS	0.647	0.314	0.972

*** Significant at the 0.01 level; ** at the 0.05 level; [◊] for *p*-values in (0.05, 0.10).

Appendix C (Hypothesis 4ii)

We report here the second part of Hypothesis 4 that was pre-registered along with the corresponding results. In our design, we have passive boss players in the *Algo* treatment so that we can evaluate their perception of distributive justice in the case of dismissal algorithms and compare it with human bosses in the *Human* treatment.

In particular, we anticipated that passive bosses would exhibit a negative reaction to algorithmic dismissals as humans generally dislike being replaced by machines (Granulo, Fuchs and Puntoni, 2019; Hidalgo et al., 2021), even though in our case we maintained bosses' revenues constant across treatments. The *Algo* treatment was a case in which the discretionary power of the boss was given to algorithms, and this task delegation to algorithms is likely to be perceived negatively (see Chugunova and Sele, 2020 for a review). Indeed, humans tend to value having control over the automated technology (Dietvorst, Simmons and Massey 2018; Gogoll and Uhl, 2018). This effect should be especially pronounced when tasks allow for subjective judgments, which is the case of bias and handicap tasks (Highhouse, 2008; Bigman and Gray, 2018) and when bosses exhibit high aversion to machines, that is a high NARS score. We state this hypothesis below.

Hypothesis 4ii.

Human bosses will perceive distributive justice to be higher in the Human than in the Algo treatment, especially for bias and handicap tasks and when their NARS score is high.

By contrast with Hypothesis 4ii, we do not find a significant interaction effect 'Algo \times NARS' (see regression [1] in Table C1) for the perceived distributive justice of bosses. However, in line with this hypothesis, we show that bosses perceived distributive justice to be lesser in the *Algo* than in the *Human* treatment (see negative and significant coefficient for 'Algo' in regression [2] in Table C1). This effect tends to be more pronounced for handicap (see 'Algo \times Handicap' in regression [2], p -value = 0.106) and bias tasks (see 'Algo \times Bias' in regression [2], p -value = 0.893) although not significantly so.

Table C1. Perceived distributive justice as a function of treatments and NARS scores.

This table presents the results of panel regressions with random effects and clustered standard errors (in parentheses) at the session level using wild bootstrapping techniques as in Cameron and Miller (2015) with 5,000 replications, which is recommended given our number of clusters (12). The Algo Dummy takes value one for the *Algo* treatment and value zero otherwise. The Bias (Unfair) Dummy takes value one when workers complete a *bias (fairness)* task. Ability is the performance on the task completed by all participants at the beginning of the experiment. We include fixed effects for task types.

Dependent Variable	Perceived distributive justice [1]	Perceived distributive justice [2]
Constant	13.441*** (3.933)	10.972*** (0.505)
Algo × NARS	-0.070 (0.271)	-
NARS	-0.124 (1.462)	-
Algo Dummy	-1.785 (1.462)	-1.604*** (0.376)
Algo × Bias Dummy	-	-0.046 (0.343)
Algo × Handicap Dummy	-	-0.571 (0.353)
Bias Dummy	-	-0.238 (0.331)
Unfair Dummy	-	-0.121 (0.276)
Period	0.004 (0.070)	-0.052 (0.052)
Ability	-0.046 (0.063)	0.015 [◊] (0.008)
N	553	553
R ²	0.123	0.103
Prob > χ^2	<0.001	<0.001

*** Significant at the 0.01 level; ** at the 0.05 level; [◊] for *p*-values in (0.05, 0.10).