



CENTRE FOR DECISION RESEARCH & EXPERIMENTAL ECONOMICS



Discussion Paper No. 2017-02

Antonio A. Arechar, Simon
Gächter & Lucas Molleman
January 2017

**Conducting interactive
experiments online**

CeDEx Discussion Paper Series
ISSN 1749 - 3293



CENTRE FOR DECISION RESEARCH & EXPERIMENTAL ECONOMICS

The Centre for Decision Research and Experimental Economics was founded in 2000, and is based in the School of Economics at the University of Nottingham.

The focus for the Centre is research into individual and strategic decision-making using a combination of theoretical and experimental methods. On the theory side, members of the Centre investigate individual choice under uncertainty, cooperative and non-cooperative game theory, as well as theories of psychology, bounded rationality and evolutionary game theory. Members of the Centre have applied experimental methods in the fields of public economics, individual choice under risk and uncertainty, strategic interaction, and the performance of auctions, markets and other economic institutions. Much of the Centre's research involves collaborative projects with researchers from other departments in the UK and overseas.

Please visit <http://www.nottingham.ac.uk/cedex> for more information about the Centre or contact

Suzanne Robey
Centre for Decision Research and Experimental Economics
School of Economics
University of Nottingham
University Park
Nottingham
NG7 2RD
Tel: +44 (0)115 95 14763
Fax: +44 (0) 115 95 14159
suzanne.robey@nottingham.ac.uk

The full list of CeDEX Discussion Papers is available at

<http://www.nottingham.ac.uk/cedex/publications/discussion-papers/index.aspx>

Conducting interactive experiments online

Antonio A. Arechar* · Simon Gächter · Lucas Molleman*

29 January 2017

Abstract Online labor markets provide new opportunities for behavioral research, but conducting economic experiments online raises important methodological challenges. This particularly holds for interactive designs. In this paper, we provide a methodological discussion of the similarities and differences between interactive experiments conducted in the laboratory and online. To this end, we conduct a repeated public goods experiment with and without punishment using samples from the laboratory and the online platform Amazon Mechanical Turk. We chose to replicate this experiment because it is long and logistically complex. It therefore provides a good case study for discussing the methodological and practical challenges of online interactive experimentation. We find that basic behavioral patterns of cooperation and punishment in the laboratory are replicable online. The most important challenge of online interactive experiments is participant dropout. We discuss measures for reducing dropout and show that, for our case study, dropouts are exogenous to the experiment. We conclude that data quality for interactive experiments via the Internet is adequate and reliable, making online interactive experimentation a valuable complement to laboratory studies.

Keywords Experimental methodology · Behavioral research · Internet experiments · Amazon Mechanical Turk · Public goods game · Punishment

JEL Classification C71 · C88 · C90 · D71

A. A. Arechar (✉)

Department of Psychology, Yale University, New Haven, Connecticut, USA

e-mail: antonio.alonso@yale.edu

S. Gächter (✉)

CeDEx and School of Economics, University of Nottingham, University Park, Nottingham NG7 2RD, UK

CESifo, Schackstrasse 4, 80539 Munich, Germany

IZA, Schaumburg-Lippe-Strasse 5-9, 53113 Bonn, Germany

e-mail: simon.gaechter@nottingham.ac.uk

L. Molleman (✉)

CeDEx, University of Nottingham, University Park, Nottingham NG7 2RD, UK

Max Planck Institute for Human Development, Lentzeallee 94, 14195 Berlin, Germany

e-mail: molleman@mpib-berlin.mpg.de

* L.M. and A.A. contributed equally to this manuscript

1 Introduction

Online labor markets such as Amazon Mechanical Turk (MTurk) are increasingly popular tools for behavioral scientists. With their large and diverse pools of people ready to promptly perform tasks for pay, these markets present researchers with new opportunities to recruit participants for experiments.¹ Studies from across the social sciences have systematically compared data collected online with data from the physical laboratory. Their conclusions are promising: classic results from psychology and economics have been replicated using online samples, and the data obtained online is deemed as reliable as that obtained via traditional methods.²

Despite its great potential, behavioral research online has so far remained largely limited to *non-interactive* decision-making tasks or one-shot games with simultaneous decisions. Current online studies of social behavior often use survey software such as *Qualtrics* or *SurveyMonkey* to document decision making in tasks that participants complete individually, and emulate interactions through post-hoc matching. Although this approach can be powerful, it does not permit the study of repeated, ‘hot’ interactions where live feedback between participants is essential. Experimental designs with live interaction are rarely implemented online, partly because there is not yet a widely-used web-based equivalent of z-Tree (Fischbacher 2007).³

In this paper we assess the potential for *interactive* online experiments, where a set of participants interacts for more than one repetition. Interactive experiments raise novel challenges throughout the whole life cycle of an experiment. Our approach is to discuss these challenges, that is, methodological differences and similarities between interactive experiments in physical and online laboratories. We discuss these step-by-step, from recruitment to dismissal of participants after the experiment.

A particularly important challenge of interactive online experiments relates to participant dropout. While in the physical laboratory participants rarely leave a session, online experiments are more prone to dropouts which affect both the participant who is dropping out and their interaction partners (who still have to be paid for their participation). If dropouts happen for reasons exogenous to the experiment – e.g. due to network problems, frozen screens, or random distractions – they are just a (costly) nuisance to the experimentalist. Much more problematic are dropouts that happen *endogenously*, that is, people quitting because of what has happened in the

¹ MTurk is by no means the only online labor market used for behavioral research. Peer et al. (2016) identified six other platforms for crowdsourcing behavioral research, in addition to platforms such as *Prolific Academic*. See Buhrmester et al. (2011), Horton et al. (2011), Rand (2012), and Paolacci and Chandler (2014) for evaluations of the potential of MTurk for conducting behavioral experiments in psychology and other social sciences.

² For example, Chesney et al. (2009) conduct classic experiments in a virtual world platform and observe behavior similar to the laboratory. Hergueux and Jacquemet (2015) find that social preferences of student participants elicited online and offline are qualitatively very similar; see their Table 1 for further references. Schmelz and Ziegelmeyer (2015) compare the behavior of student participants in a principal-agent experiment conducted in the laboratory and on the Internet and conclude that data can be of similar quality. On MTurk, Horton et al. (2011) replicate classic framing effects and report similar levels of cooperation in a one-shot prisoner’s dilemma played online and in the laboratory. Berinsky et al. (2012) replicate classic findings on framing and risk tasks. For overviews of replication studies, see Paolacci et al. (2010), Behrend et al. (2011), Klein et al. (2014), and Mullinix et al. (2015).

³ Thus far interactive research developed its own bespoke software (see: Egas and Riedl 2008; Suri and Watts 2011; Wang et al. 2012; Gallo and Yan 2015; Nishi et al. 2015; Schmelz and Ziegelmeyer 2015; Stagnaro et al. 2016), but recent software is likely to change this trend (see: BreadBoard, MobLab, oTree, SoPHIE, and UbiquityLab).

experiment. Such dropouts could jeopardize the internal validity of experiments (Zhou and Fishbach 2016).

As a case study we replicate a repeated public goods game with and without peer punishment used in cross-cultural research (Herrmann et al. 2008), utilizing a sample of US participants recruited via MTurk.⁴ We chose to replicate this experiment because it is fairly long and logistically complex. It is a within-subjects design with two experimental conditions of ten periods each, where, after the first set of ten periods, participants receive new instructions. Moreover, this experiment has often been replicated, and its design allows us to evaluate whether dropouts depend on the experimental conditions (that is, the presence or absence of punishment). We report data from participants recruited via MTurk (62 groups) and participants from the physical laboratory (18 groups). We used our own software LIONESS (Sect. 2.5), developed for conducting interactive online experiments.

We observe that basic patterns of behavior online are similar to those in the laboratory. In the absence of punishment, aggregate levels of cooperation are higher on MTurk than in the laboratory, but show similar rates of decay over time. Moreover, our econometric analysis reveals that in both of our samples the group contributions strongly determine the level of cooperation. The introduction of punishment promotes the emergence and maintenance of cooperation in both samples. Punishment is mainly prosocial in nature in both samples (cooperators punish non-cooperators) but occurs less frequently online.

Our most important result is that, in our implementation, dropouts are most likely due to reasons that are exogenous to the experiment. Together with the replication of findings from the laboratory, our results suggest that online interactive experiments can be a reliable tool for collecting internally-valid data and hence are a valuable complement to the physical laboratory.

Our paper contributes to a recently-emerged literature on the reliability of data gathered on online labor markets such as MTurk (see references in footnotes 1-3). The most important predecessor of our paper is Anderhub et al. (2001), who compared online and laboratory experiments in the very early days of experimentation on the Internet. They also provide a methodological discussion that, however, could not take into account the specific properties of modern online labor markets where the bulk of present-day online experimentation is happening.⁵

The remainder of this paper is structured as follows. In Sect. 2, we introduce the experimental design. In Sect. 3, we discuss the conceptual and logistical differences between conducting interactive experiments in the laboratory and online and outline our approach for dealing with them, highlighting important aspects of the data-collection process relating to attention and attrition. Sect. 4 shows the results of our experiment, systematically comparing cooperation and punishment behavior in our two samples. In Sect. 5 we present a detailed analysis of attrition in our online experiment. Finally, in Sect. 6 we make concluding remarks.

⁴ In one of the first studies using economic games on MTurk, Suri and Watts (2011) replicate the laboratory results found by Fehr and Gächter (2000) using a repeated public goods game without punishment.

⁵ Another early paper on experimentation on the internet is Reips (2000). The author discusses the advantages and disadvantages of conducting online experiments and also provides a methodological discussion of how internet experiments can be used to validate laboratory data.

2 A case study to compare online and laboratory experiments

We base our discussion of online and laboratory experiments on a well-established paradigm: a public goods game with and without punishment (Fehr and Gächter 2000; Fehr and Gächter 2002). In this section we present the design of our experiments conducted in the laboratory and replicated online with a sample of participants recruited via MTurk. For the laboratory and the online samples, instructions and experimental screens were identical (screenshots are presented in Appendix B).

2.1 General setup

Our experiment follows the within-subject design of Herrmann et al. (2008) and implements a repeated four-person public goods game with two conditions: one without punishment followed by one with punishment. Groups were constant throughout the experiment ('partner matching'), and each condition ran for ten periods. Participants were aware that there were two 'parts' to the session (which corresponded to the conditions without and with punishment) but learned about the details of the second part only after the first one had finished.

At the beginning of a session, participants read on-screen instructions for the first experimental condition: the public goods game without punishment. Experimental instructions were shorter than those in Herrmann et al. (2008) (see Sect. 3 for rationale). Participants could start the interaction phase only once they had completed a set of comprehension questions.

2.2 Condition 1: A public goods game without peer punishment

In each period of the 10-period game, all four group members received an endowment of 20 points and simultaneously decided how many of those points to keep for themselves, and how many to contribute to a 'group project' (i.e. the public good). After all members had made their decision, the sum of all contributions was multiplied by 1.6 and distributed equally among all group members irrespective of their contributions. This setup reflects a social dilemma: in a given period overall earnings are highest when each of the group members contributes all 20 of their points to the public good, while individuals maximize their earnings by contributing 0 regardless of the contributions of the others. Once all contributions had been made participants learned the result of that period. Apart from their own contribution and earnings, they were informed of the average contribution in their group. Subsequently, a separate screen showed the contributions of each of their fellow group members.

2.3 Condition 2: A public goods game with peer punishment

Once the 10 periods of Condition 1 were over, participants received new on-screen instructions about Condition 2. This condition also consisted of ten periods and was completed in the same groups as Condition 1. Again, the periods started once all group members had completed the comprehension questions. The decision situation was similar to Condition 1, but we introduced one change: once participants learned the contributions of each of their group members, they were allowed to assign up to 10 deduction points to each of their peers. Each assigned deduction point resulted in a loss of 1 point for the participant assigning it, and a loss of 3 points for its

target. At the end of each period a separate screen informed participants of the total number of points they assigned and received. In cases where a participant made a loss during a period, only the costs of assigning deduction points would count towards the final earnings (cf. Herrmann et al. 2008). Each session concluded with a questionnaire including demographic items.

2.4 Online and laboratory sample

In all sessions, participants received instructions and made their decisions via web browsers. The program was implemented in the experimental software LIONESS (Sect. 2.5). Both online and in the laboratory, sessions took 28 minutes on average. This is considerably shorter than the original study by Herrmann et al. (2008), but longer than typical tasks on MTurk.

For our *online* sample we recruited participants via MTurk, restricting their geographical location to the USA (for comparability with our laboratory sample, see below). Results are based on 24 sessions, with 248 participants in total (62 groups of four). The average age of participants in this sample was 31.5 years (s.d. 9.06), and 38.1% were female. Average earnings in our online sample were \$6.69 (s.d. \$1.03), which were paid via MTurk.

The data from our *laboratory* sample were collected at universities in two different cities in the USA (Harvard University, Boston MA; and Yale University, New Haven CT) over 8 sessions, with 72 participants in total (18 groups). Laboratory participants were invited through e-mails using the online recruitment software SONA. The average age for participants in this sample was 25.2 years (s.d. 7.45) and 41.7% were female. Average earnings in our laboratory sample were \$20.02 (s.d. \$1.65), paid in cash upon session completion. In order to conform to standards of the respective laboratories and average expected wages on MTurk, we used an exchange rate of \$0.02 in the laboratory and \$0.01 in our online sessions; show-up fees were \$10 and \$1 in the laboratory and online, respectively.⁶

2.5 The software used to conduct interactive experiments: LIONESS

We conducted both the laboratory and online experiments with LIONESS (Live Interactive Online Experimental Server Software). LIONESS provides a basic architecture for conducting interactive experiments online. Its key features reflect the solutions to the logistical challenges discussed in Sect. 3; dynamically grouping participants to minimize waiting times, regulating interactions in groups, promoting participants' attention to the experiment and dealing with participants dropping out of an experiment. The software developed for the experimental conditions reported here can be downloaded at: <http://lioness.nottingham.ac.uk>.

⁶ Using four canonical one-shot games, Amir et al. (2012) find results comparable to those found in the laboratory, even when using relatively low stakes. In the public goods context, Suri and Watts (2011) also report similar levels of cooperation when using two different compensation levels.

3 Methodological differences in conducting interactive experiments in the laboratory and online

3.1 The online laboratory MTurk

While our discussion of online experiments is based on an MTurk sample, many issues also hold for other online platforms (cf. footnote 1). MTurk is a very large online labor market, which offers an active pool of over 500,000 workers. The MTurk workforce completes over 40,000 Human Intelligence Tasks (HITs) every day (www.mturk-tracker.com; Difallah et al. 2015; Ipeirotis 2010). MTurk ‘workers’ browse HITs that are published by ‘requesters’ who provide a brief description of the task, its expected duration and the minimum payment workers will receive upon completion (see Appendix B for screenshots of the HIT as published on MTurk). HITs typically involve short individual assignments which computers are currently unable to perform (Berinsky et al. 2012), such as the processing of images or data cleaning. Due to the sheer size of the pool of workers ready to perform tasks for pay, MTurk enables researchers to conduct large-scale experiments and to implement an effective random assignment of participants to different conditions beyond the capacity of a typical physical laboratory.⁷ While not primarily designed for academic research, MTurk has the potential for efficient data collection. As mentioned in the introduction, researchers have reported promising results from questionnaire studies and experiments without repetitions (e.g. one-shot Prisoner’s dilemma), with results from MTurk comparable to those obtained from laboratory samples (e.g. Paolacci et al. 2010; Horton et al. 2011; Goodman et al. 2013).

Despite their promise, online behavioral experiments have conceptual and logistical challenges that are usually not present in the laboratory. Here we focus on differences between laboratory and online experimentation that are specific to *interactive* designs. See Buhrmester et al. (2011) and Paolacci and Chandler (2014) for extensive discussions of differences regarding non-interactive (survey-style) designs.

3.2 A step-by-step comparison of laboratory and online experiments

We now discuss the implementation of interactive experiments in the laboratory and online via MTurk, based on the design we presented above. Our discussion is chronological in the way a typical experiment proceeds from recruitment to dismissal of participants. The four following subsections discuss the main phases of a typical experimental session (recruitment, session start-up, interactive decision making and payment). Along the way we highlight the extent to which our approach bridges these gaps in the (relatively long and logistically-challenging) experiment presented above. Table 1 provides a concise overview of the issues discussed in this section.

3.2.1 Recruitment

In a typical laboratory experiment, participants receive a show-up fee for attending. Still, the main part of the participant’s payment is usually determined by the decisions made over the

⁷ In comparison to the student or community samples normally used in the laboratory, Mturk samples are also more diverse (Krupnikov and Levine 2014; Paolacci and Chandler 2014; Weinberg et al. 2014; Berinsky et al. 2012).

course of the session. In a typical task on MTurk, participants are paid a flat reward per HIT, and the part of the earnings determined by their decisions can be added to their payments as a ‘bonus’. Consequently, a HIT that pays a relatively large flat fee usually draws more attention than one that promises a large bonus. This is particularly relevant for interactive experiments where participants need to wait for others to form a group at the start of a session (see below).

Phase / challenge	Laboratory	Online (MTurk)
<i>Recruitment</i>		
Show-up fees	Typically a small part of total payoffs. Guaranteed when participant shows up to the session.	Relatively large show-up fees promote recruitment rates, thereby facilitating prompt group formation. Experimenter can approve or reject the task submitted; if rejected no fee is paid.
Inviting participants	Invitations sent well in advance, participants commit to a session. Recruitment often from a pre-existing database.	Sessions advertised online as HITs and can be completed immediately.
Selection into the experiment	At sign-up, participants know very little about the experiment. Details of the task are communicated once participants are in the laboratory.	Experiments are typically advertised as HITs with a brief task description. ‘Workers’ browse available HITs and accept those of their preference.
Experienced participants	Invitation conditioned on well-defined criteria of the laboratory’s records.	HITs targeted at subsets of MTurk workers; experimenter can specify exclusion criteria. Many MTurk workers will have participated in many prior studies.
<i>Session start-up</i>		
Duplicate participants	Registration protocols usually prevent duplicate participation.	Amazon acts against multiple worker accounts, but they exist.
Comprehension	Participants can ask questions; comprehension questions ensure understanding.	Experimenter is physically absent and can’t answer questions directly. Compulsory comprehension questions can be added but may make experiment (too) long for some participants.
<i>Experimental interactions</i>		
Forming groups	Easy to guess how many participants will attend; group settings can be pre-defined.	Hard to guess how many participants will attend; groups can be constructed ‘on the fly’.
Deception	In experimental economics deception is prohibited and laboratories foster reputations for non-deception.	Because all requesters use the same subject pool, some participants may have experienced deception because requesters from other disciplines may use it.
Communication	Hardly an issue; experimenter can restrict communication between subjects.	Participants may in principle collude through external channels though this is difficult in practice.
Experimental flow	Closed form software like z-Tree specifies session progress.	Scripted browser navigation specifies progress.
Attrition (‘dropout’)	Hardly an issue; participants that start a session usually finish it.	Major challenge to internal validity, if dropout rates vary with treatment, selection bias may arise.
<i>Payment</i>		
Payments	Cash usually paid upon completion.	Automatic transfer through Amazon.
Cost per participant	Relatively high but predictable.	Relatively low but varies with attrition.

Table 1 Methodological differences in conducting interactive experiments in the laboratory and on MTurk

Sessions in the laboratory are pre-scheduled. A database contains the contact details of a pool of aspiring participants, who can register (and cancel) within a determined time window. Pre-scheduling ensures that the number of participants can be anticipated quite accurately before a session takes place, and including a small number of backup participants can prevent problems associated with unannounced non-attendance. Online platforms such as MTurk allow for instant recruitment of participants, facilitating a time-efficient method of data collection.⁸

To take advantage of these opportunities, LIONESS was ready to accommodate new participants during a time window specified by us, while capping the maximum number of entrants. In our experiment, we invited participants to sign up within 20 minutes of the HIT being posted and allotted them 45 minutes to complete the task. In addition, we asked them to start immediately and did not let anyone enter once the number of participants reached a limit. Recruitment rates were high (in a typical session with 100 slots, the first 50 participants typically entered within the first 5 minutes after the publication of the HIT), facilitating prompt group formation once participants had read and understood the instructions.⁹

Invitations to laboratory experiments do not typically reveal any information on the contents of the experiment. On MTurk, however, participants browse various tasks that are currently available to complete for pay. This requires a HIT description giving the workers some idea of what the task will involve. To avoid self-selection (based on the topic of the experiment) into interactive experiments as much as possible, an experimenter can leave out any detailed information in the HIT description.¹⁰ For example, we did not announce that these are public goods experiments. On the other hand, it is essential that workers know that the HIT will involve live interactions with other people, and as a consequence they are expected to complete the interactive HIT without delay and without interruptions. The HIT description within MTurk is an appropriate place to make participants aware of this (see Appendix B for screenshots).

Typical laboratory subject pools are replenished annually when a new cohort of first-year students arrives on campus, and recruitment software allows invitation of only those participants with no (or little) experience with the experimental paradigm of a study. By contrast, many MTurkers are replenished more regularly, but quickly acquire a sizable experience participating in hundreds of academic studies of all kinds.¹¹ While researchers may have reason to believe that (frequent) prior experience could be an issue for their experiment (for a discussion of “lab rats” see Guillen and Veszteg 2012), MTurk facilitates inviting participants based on various criteria (e.g. number of HITs completed, their success rate or their geographic location). In addition,

⁸ Also, the sheer size of the potential pool of participants on MTurk allows researchers to use designs requiring large groups, and to obtain sufficient statistical power to detect effects of relatively small size (e.g. Hauser et al. 2016).

⁹ We chose a 20-minute time slot based on pilot sessions. Over time, MTurk HITs move down the list of most recently-published jobs and entrance rates tend to taper off. When participants trickle in at low rates, the chances that they have to wait for a long time before they can be matched into a group increase substantially.

¹⁰ Self-selection into experiments may not be a specific problem for online studies. Selection effects may also occur in physical laboratories which regularly run interactive tasks. Participants may anticipate their tasks involving interactions with other participants in their session due to previous experience or hearing about it from other members of the subject pool; see Krawczyk (2011), Anderson et al. (2013), Cleave et al. (2013) and Abeler and Nosenzo (2015) for detailed discussions.

¹¹ (Stewart et al. 2015) estimate a worker half-life of about seven months and Rand et al. (2014) report a median of 300 academic studies, 20 of which occurred in the past week.

post-experimental questionnaires can include self-reported measures of participants' familiarity with decision-making experiments and specific experimental paradigms.

In our case, we used MTurk's options to restrict the geographical location of the participants to the US for comparability with our laboratory sample. In addition, to increase the likelihood that participants completed our HIT with care, we only allowed workers with at least 90% of their previous HITs approved by requesters to participate (see Peer et al. 2014 for a detailed discussion of approval rates).

3.2.2 Session start-up

For many studies it is essential that participants only take part once. In the laboratory it is relatively straightforward to implement this, particularly if the experimenter uses recruitment software such as *ORSEE* (Greiner 2015), *SONA* or *hroot* (Bock et al. 2014) and is physically present during laboratory sessions to confirm identities. For online sessions, however, re-takers may seriously compromise the data (e.g. by operating two browsers within the same experiment, potentially even controlling two players within the same group). Accordingly detecting them requires specific measures. Within a session, we prevented duplicate participation by logging the user's IP address and blocking users that had already been connected to the experimental server. Between sessions, we used third-party software to prevent workers who had already participated in a specific HIT from being invited for future sessions.¹²

In a typical laboratory session participants are allowed to ask questions which the experimenter can answer in private, before the interactive phase of the experiment begins. For online sessions this is not feasible. To ensure that participants had a thorough understanding of the experimental decision situation and did not rush through the instructions, we introduced compulsory comprehension questions which participants had to solve before entering the decision-making phase of the experiment.¹³ A fraction of participants who entered the experimental pages did not proceed past the instructions and never reached the comprehension questions. In our online sample, 83.2% of the individuals who did reach the comprehension questions solved them successfully.¹⁴

3.2.3 Experimental interaction

In the laboratory, all participants typically arrive at a session at the same time and will simultaneously complete comprehension questions. By contrast, participants in online sessions may arrive during a time window set by the experimenter (20 minutes in our case), and the timing of completing comprehension questions may therefore vary substantially. Accordingly,

¹² We used *UniqueTurker* (<http://uniqueturker.myleott.com>) to prevent re-takes between sessions. This method uses a unique MTurk identification number which is linked to each worker and that Amazon constantly monitors to avoid duplicate participation. Alternatives for these methods include *Turkprime* (<http://turkprime.com>).

¹³ We used considerably shorter experimental instructions than Herrmann et al. (2008) to minimize the variation in the time that participants take to read them; with dynamic group formation, minimal variation in preparation time is highly desirable (e.g. to reduce attrition, see below).

¹⁴ When entering an incorrect answer, subjects were allowed to try again. Incidentally, the total number of incorrect attempts in the six comprehension questions prior to Condition 1 did not differ between our laboratory and online samples (4.68 vs 4.46, respectively, Mann-Whitney test: $p=0.245$).

we formed groups ‘on the fly’: participants who successfully completed the comprehension questions waited in a ‘lobby’. As soon as this lobby contained sufficient participants, a group was formed and its members were sent to the interaction phase of the experiment. An alternative (yet considerably less time-efficient) approach is to run pre-tests with participants and to build a ‘standing panel’ from which candidates for experimental sessions are recruited (see Suri and Watts 2011; Gallo and Yan 2015).

Although deception is uncommon in experimental economics, participants on MTurk are likely to encounter studies using deception (e.g. Pfattheicher and Schindler 2015). Participants may therefore be skeptical about the truthfulness of experimental instructions and doubt if their interaction partners are real people and not robots pre-programmed by the experimenter publishing the HIT. To promote trust between us (as experimenters) and the participants, we continuously strive to maintain a good reputation on our MTurk requester account (our records and those of other requesters can be found at <https://turkopticon.ucsd.edu>). In addition, our HIT description stated explicitly that groups were formed of real people recruited from MTurk.

In order to keep the attention of the participants focused on the experiment (and not have them dropping out in the very first period of the game), we clearly communicated the number of other participants they were waiting for at any given moment, and we added an on-screen countdown indicating the maximum amount of time left before participants could choose to leave the experiment if no group could be formed. When this timer reached zero in our experiment, participants could choose to either return to the lobby and wait for two additional minutes or to leave the session and collect their participation fee (of \$1). This procedure led to a total of 89% of participants who correctly completed the comprehension questions being successfully matched into a group and starting the interaction phase.¹⁵ The remaining 11% could not be matched in a group of four, and were paid their participation fee.

While in the laboratory the experimenter can monitor and enforce any restriction of communication between participants, it is in principle harder to categorically exclude the possibility that online participants communicate with their interaction partners through external channels. We ran relatively large sessions in which participants could not be identified to prevent them from colluding via online forums such as *Reddit* or *MTurk Crowd*. In fact, forum discussions, which are usually moderated and prohibit the dissemination of the content of HITs and the discussion of strategies, typically center upon the attractiveness of a HIT in terms of earnings and length rather than its content (Arechar et al. 2015; Chandler et al. 2014). Therefore, communication between participants is a potential problem for online experiments, but it is not any more severe for interactive designs. A similar argument could be made for communication with other people who are not participating in the task at hand (e.g. someone else in the room while completing the task). Although such communication is harder to control in online experiments than in the physical laboratory, this issue is not specific to interactive designs either.

¹⁵ As subjects completed the experiment in their browsers, we strictly regulated navigation between experimental pages. Participants navigating the pages at will may lead to serious disruptions in the experimental flow and even lead to crashing sessions. We used LIONESS to regulate the sequence of an experiment in a browser by using a pop-up window without navigation bars and overwriting the browser history within that window when moving on to another page. This procedure also prevented participants from revising their decisions during a given interaction.

The most severe problem for online interactive studies, and the largest discrepancy with laboratory experiments, is *attrition* (participant dropout). In laboratory sessions participants very rarely leave or turn out to be unable to complete a session. In online experiments, attrition is a major issue – there is no straightforward way to prevent participants from leaving a session by closing their browser window or failing to submit responses to experimental decision situations due to technical problems. Moreover, in contrast to the laboratory, interaction partners are geographically scattered and the progress of an experiment depends on their joint attention to it. Typically, groups proceed at the pace of the slowest participant and long waiting times increase the risk of reduced attention, which may ‘cascade’ through the group. Thus, we took measures to retain attention and promote successful completion.¹⁶ In our case, we used on-screen timers and told participants that failure to reach a decision in due time would result in their removal from the experiment without payment. Furthermore, in the event of a group member dropping out we notified the remaining participants of that and they continued in their reduced group.¹⁷ Data from incomplete groups is not included in the results reported in Sect. 4. An alternative approach to dealing with dropouts is to terminate the whole group once a member drops out. This may, however, damage the reputation of the experimenter as participants will be unable to earn as much as they had anticipated.

Our procedure of letting smaller groups continue ensures that real people generate all the information that participants respond to. Alternative solutions to non-responding participants, such as introducing random decisions or repeating previous decisions (e.g. Suri and Watts 2011; Wang et al. 2012), may affect the behavior of those who are still in the experiment (now responding to partially computer-generated information) which potentially compromises the internal validity of the data from groups affected by a dropout. It also raises issues of deception if such computer-generated information is not disclosed. Moreover, this procedure may also compromise the validity of data from groups unaffected by attrition, as participants cannot know whether their interaction partners’ behavior shown to them is real or generated by a computer.

Due to the nature of conducting research via the Internet, some level of attrition seems unavoidable. Attrition rates are likely to vary with factors such as group size, complexity of the decision situation, and the pace of the experiment.¹⁸ Despite our measures to prevent attrition, 84 participants (18%) who started the interaction phase dropped out at some point in our

¹⁶ On waiting pages (to which LIONESS directed participants once they made their decisions and waited for their fellow group members to do so as well), we added an animation (a commonly used circular icon spinning around) to assure the participants that the experiment was still active. Additional measures to promote attention may include playing sounds (e.g. Hauser et al. 2016) or having the browser window with the experimental pages placed in the foreground of the participants’ computer screen (‘focus stealing’) when a decision is due. We did not use these more intrusive measures in this study.

¹⁷ The experimental instructions mentioned that the total number of points contributed to the public good was multiplied by 1.6 before distributing the resulting amount equally among the group members. A group member dropping out therefore increased the marginal per capita return of contributions.

¹⁸ Note that both the stability of the experimental software and the stability of the server on which the software is run might strongly affect attrition too. Pilot sessions led us to deploy LIONESS on high-CPU servers (available through Google Cloud) to facilitate many simultaneous server-client interactions.

experiment.¹⁹ As these participants were distributed across experimental groups, the fraction of the data set affected by these dropouts was considerably larger.

Figure 1 tracks the distribution of group sizes over time. All groups are initially formed of four group members, but a group’s size may decrease over the course of the experiment as participants drop out. The figure shows that our experiment suffered from quite substantial attrition and only 53% of the groups finished with all four members. Loss of group members was particularly likely around the ‘waiting room’ stages preceding periods 1 and 11. Specifically, in period 1 dropouts are presumably increased due to participants losing attention while waiting for their group to form. Similarly, before period 11 started participants had to wait until each of their group members had completed the comprehension questions, which could take a considerable amount of time. In some cases this led to the termination of the whole group.

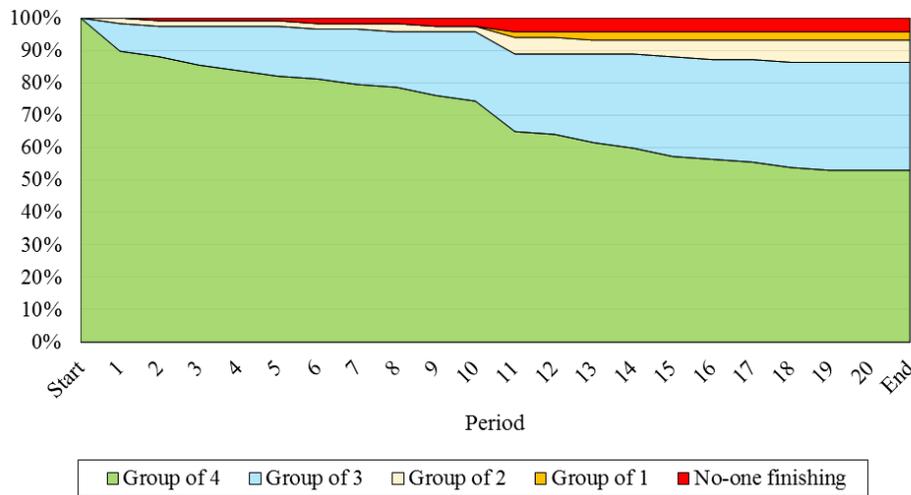


Fig. 1 Attrition throughout the course of the experiment. Colors depict the group size. We always started with groups of four but let participants continue if a member dropped out.

3.2.4 Payment and costs of experiments

After a session is over, participants are typically paid according to their performance. One of the advantages of MTurk is that it facilitates secure payments. The experimental software can generate a random code for each participant which can be matched with the participants’ MTurk ID, allowing for payments according to performance in the experiment. It is important to process payments immediately in order to maintain a good reputation as an MTurk requester.

All in all, typical costs per useable data point in an experiment with participants recruited via MTurk are likely to be lower than in a laboratory experiment. Nevertheless, these costs may vary with attrition rates, which can be affected by the specific features of the experiment such as its length and group size (as one dropout may compromise the data of the whole group). In our case, laboratory participants earned \$20.02 on average. Therefore, with a group of four as the unit of observation, a useable data point cost us \$80.08. Corresponding costs online were \$49.25.²⁰

¹⁹ Participants that did not respond in time or dropped out for any other reason did not receive any payment.

²⁰ We incurred additional costs in our online sessions through: (i) 69 participants who did complete the comprehension questions but were not allocated to a group and who received \$1; (ii) participants from groups with

4 Results

4.1 Contribution behavior

Figure 2 shows the aggregate contribution dynamics for both the online and laboratory samples. In the condition without punishment, overall contributions were higher in our online sample than in the laboratory (12.52 vs 8.30, $p=0.003$).²¹ This result is consistent with recent literature finding college students to be less cooperative than non-students and other adults (Belot et al. 2010; Carpenter and Seki 2011; Gächter et al. 2004; Stoop et al. 2012). The difference in contributions emerges right in the very first period of the game, with online participants contributing substantially more to the public good (15.00 vs 11.04, $p=0.001$). Higher contributions by MTurkers can be only partially explained by the higher average age in the MTurk sample (OLS fitted to contribution decisions in the first period of the first condition: age $\beta=0.104$, $p=0.031$; MTurk dummy: $\beta=3.110$, $p=0.004$).²²

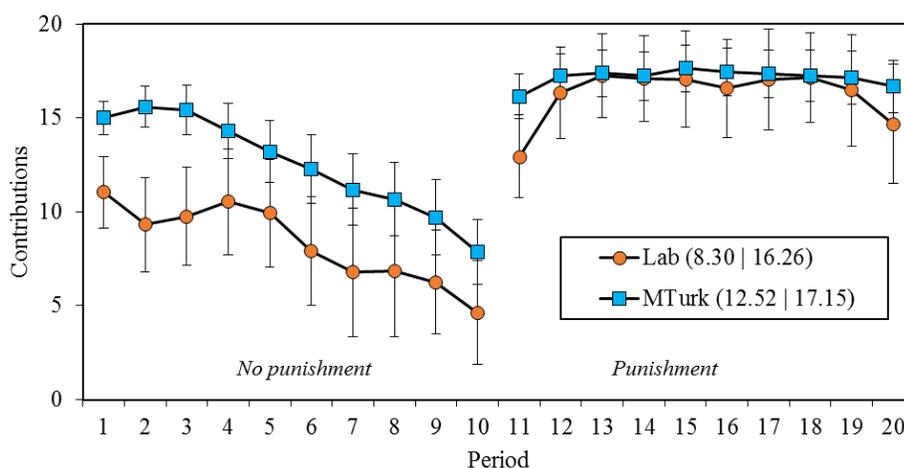


Fig. 2 Contributions over time. Numbers in parentheses are the mean contributions in each experimental condition. Error bars indicate 95% confidence intervals (clustered at the group level).

The introduction of punishment opportunities strongly increases average contributions in both samples (average contributions in periods 10 vs 11: laboratory: 4.63 vs 12.94, $p=0.001$; online: 7.85 vs 16.15, $p=0.001$). Moreover, average cooperation levels over the course of the game are higher than in the absence of punishment (average group contributions in periods 11-20 vs 1-10: laboratory: 16.26 vs 8.30, $p=0.001$; online: 17.15 vs 12.52, $p=0.001$). As in the condition without punishment, overall contributions within groups were slightly yet significantly higher in our online sample than in the laboratory (16.26 vs 17.15, $p=0.008$).²³

dropouts (a total of \$885.05); (iii) Amazon MTurk fees (\$508.92). The online costs per usable data point are computed as: (payments to the participants in the 62 groups finishing without dropouts + additional costs from *i-iii*) / 62. This amount does not include payments to subjects who took part in pilots we ran to determine our procedures.

²¹ Unless otherwise stated, all statistical tests are two-sided Mann-Whitney tests with group averages over all periods as independent observations.

²² See List (2004) for a detailed discussion of age effects on cooperation behavior. These results are also consistent with Gächter and Herrmann (2011) who found in a laboratory one-shot public goods game that older participants contribute more to the public good than younger ones.

²³ See Appendix C for a comparison of our samples with the cross-cultural samples from Herrmann et al. (2008).

The regression models in Table 2 confirm that Mturk participants initially contribute more to the public good than their laboratory counterparts (Wald test on ‘Constant’: $p < 0.001$; Table 2, columns 1 and 2). Over the course of the game cooperation decays at comparable rates (Wald test on ‘Period’: $p = 0.624$; columns 1 and 2). In the punishment condition, the constant differs between the two samples (Wald test on ‘Constant’: $p = 0.084$; Table 2, columns 4 and 5), but the effect of period is equivalent (Wald test on ‘Period’: $p = 0.407$; Table 2, columns 4 and 5).²⁴

<i>Contributions to the public good</i>						
	No punishment			Punishment		
	Laboratory	MTurk	Pooled	Laboratory	MTurk	Pooled
Period	-0.900*** (0.309)	-1.074*** (0.187)	-1.037*** (0.160)	1.139 (0.710)	0.514* (0.289)	0.682** (0.282)
Final period	-3.400 (2.253)	-2.292** (0.958)	-2.512*** (0.881)	-10.203** (4.881)	-4.184** (1.688)	-5.795*** (1.797)
MTurk			5.421*** (1.867)			4.193 (4.904)
Constant	10.470*** (1.592)	17.046*** (0.624)	11.402*** (1.650)	25.980*** (3.898)	35.272*** (3.792)	29.601*** (4.232)
N	720	2480	3200	720	2480	3200
F	8.75	33.66	34.45	2.19	3.12	3.75

Table 2 Cooperation dynamics. Tobit estimation with left-censoring for ‘No punishment’ and right-censoring for ‘Punishment’. ‘Period’ is period number; ‘Final period’ is a dummy for last period; ‘MTurk’ is a dummy for the MTurk sample. Robust standard errors clustered on groups; * $p < 0.1$, ** $p < 0.05$, *** $p < 0.01$.

Individual responses to the contributions of fellow group members were significant and similar across the online and laboratory samples (Wald test on ‘Mean peer contribution in $t-1$ ’: $p = 0.505$; Table 3, columns 1 and 2), which suggests that individual decision making online was not more ‘random’ than in the laboratory.²⁵

<i>Contribution to the public good (no punishment)</i>			
	Laboratory	MTurk	Pooled
Period	-0.401** (0.204)	-0.503*** (0.094)	-0.485*** (0.085)
Final period	-2.826 (1.941)	-1.316 (0.827)	-1.600** (0.757)
Mean peer contribution in $t-1$	0.953*** (0.125)	1.043*** (0.060)	1.027*** (0.054)
MTurk			0.759 (0.778)
Constant	-0.830 (1.674)	-0.006 (1.237)	-0.696 (0.931)
N	648	2232	2880
F	29.05	163.74	177.16

Table 3 Cooperation dynamics (no punishment). Left-censored Tobit estimation. ‘Period’ is period number; ‘Mean peer contribution in $t-1$ ’ is the average contribution of the other members in the group in $t-1$; ‘MTurk’ is a dummy for the MTurk sample. Robust standard errors clustered on groups; * $p < 0.1$, ** $p < 0.05$, *** $p < 0.01$.

²⁴ An additional regression model including an interaction between MTurk and period is not significant in either condition ($p > 0.560$). Our results are also robust to the inclusion of demographic controls and the use of a different (multilevel mixed effects) model specification (see Table A1 of Appendix A).

²⁵ An additional model with an interaction between ‘MTurk’ and ‘mean contributions’ shows no significance for this item ($p = 0.294$). As in Table 2, results are robust to different model specifications (see Table A2 of Appendix A).

4.2 Punishment behavior

Participants in the online sample punished less often than their laboratory counterparts (overall punishment frequency: 0.072 vs 0.167, $p=0.001$).²⁶ Moreover, Figure 3 shows that in both samples the frequency of punishment tends to decrease over the course of the game, albeit less markedly in our online sample. Accordingly, mean efficiency in the experimental condition with punishment was higher on MTurk than in the laboratory (averages 26.91 vs 22.64 points per individual per period, $p=0.002$). In cases where participants decided to punish, they did so equally severely in both samples. The average number of points assigned when punishing did not differ significantly between MTurk and the laboratory (4.15 vs 3.88; $p=0.545$).

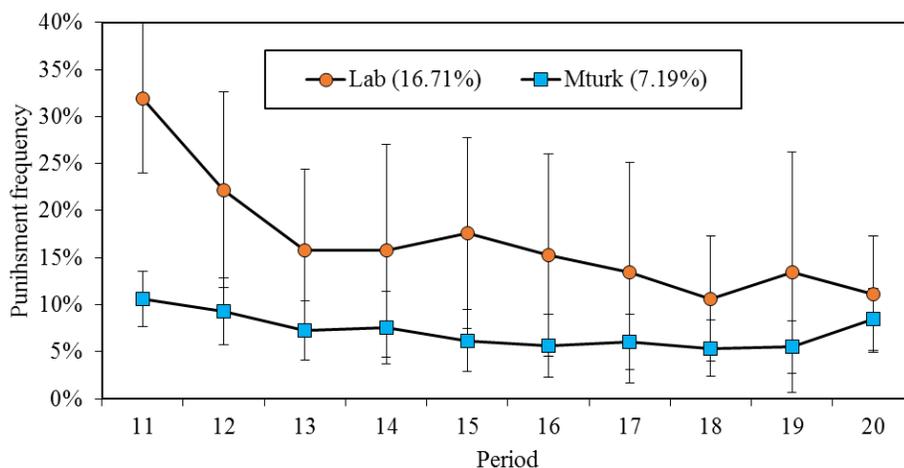


Fig. 3 Frequencies of punishment over time. Frequencies are calculated by counting instances of assigning non-zero deduction points out of the total number of punishment opportunities per participant, per recipient, per period. Mean punishment frequencies in parenthesis. Error bars indicate 95% confidence intervals clustered on groups.

Figure 4 reveals that, in both samples, punishment was predominantly pro-social in nature: most instances of punishment represented cooperators (who contributed relatively many points to the public good) punishing defectors (who contributed relatively fewer points). In the laboratory and online, both frequency and severity of punishment were higher with increasing differences in contributions between the punisher and his target (Fig. 4, compare the bottom four stacked bars; see Table A3 of Appendix A for regression analyses). Interestingly, we observe some instances of anti-social punishment in both samples (Fig. 4, top four stacked bars).

Table 4 presents an econometric analysis of punishment behavior. It confirms the observations from Figures 3 and 4. The overall patterns of punishment are similar in the laboratory and online. This analysis further reveals that online participants tended to punish less frequently and less severely, even after controlling for relative contributions and previously received punishment (Table 4, Wald test on constants; columns 1 and 2: $p=0.014$; columns 3 and 4: $p=0.064$).²⁷

²⁶ The difference remains significant in a logit model fitted to decisions to punish (no/yes) that controls for period, participant's contribution, and punishment received in the previous period (MTurk dummy: $\beta=-0.898$, $p=0.001$).

²⁷ In Table A3 of the Appendix we present a more detailed analysis of pro-social and anti-social punishment, indicating that by-and-large, punishment is driven by the same factors in the laboratory and online. These models

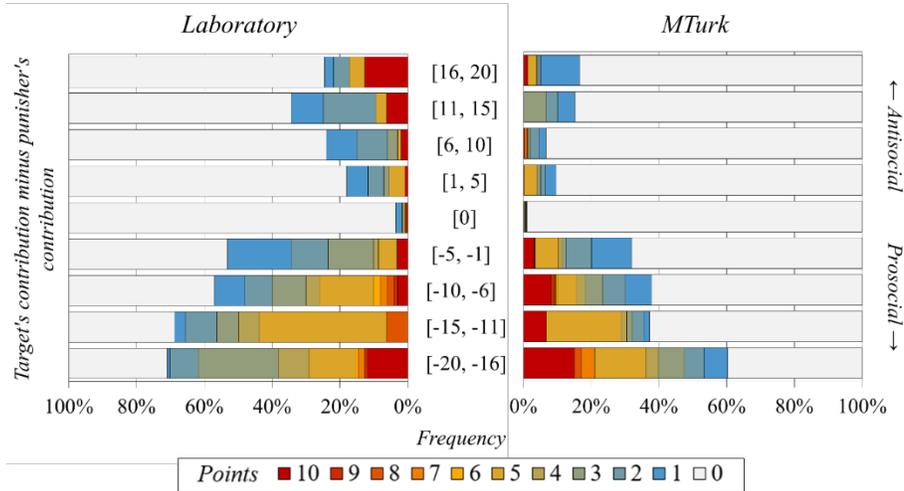


Fig. 4 Directionality and severity of punishment in our laboratory and online samples. Stacked bars show frequency distributions of punishment decisions. Each bar shows the distribution for a given difference between punishers and their target’s contribution to the public good.

	<i>Decision to punish (0=no; 1=yes)</i>			<i>Punishment severity</i>		
	Laboratory	MTurk	Pooled	Laboratory	MTurk	Pooled
Target’s contribution	-0.181*** (0.034)	-0.216*** (0.015)	-0.203*** (0.016)	-0.505*** (0.074)	-0.717*** (0.051)	-0.641*** (0.043)
Punisher’s contribution	-0.014 (0.034)	0.003 (0.026)	-0.002 (0.022)	-0.067 (0.106)	-0.011 (0.074)	-0.032 (0.063)
Mean contrib. others	0.040 (0.028)	0.065** (0.025)	0.058*** (0.019)	0.136 (0.094)	0.228*** (0.071)	0.197*** (0.058)
Rec. punishment in <i>t</i> -1	0.090*** (0.022)	0.097** (0.045)	0.092*** (0.026)	0.310*** (0.074)	0.273** (0.128)	0.284*** (0.069)
Period	-0.126** (0.056)	-0.102*** (0.037)	-0.111*** (0.031)	-0.311** (0.143)	-0.280** (0.110)	-0.289*** (0.088)
Final period	-0.633* (0.336)	0.524* (0.294)	0.150 (0.266)	-1.014 (1.220)	2.401*** (0.888)	1.271 (0.780)
MTurk			-0.965*** (0.231)			-2.631*** (0.803)
Constant	0.960* (0.496)	-0.440 (0.290)	0.671** (0.310)	1.569 (1.533)	-1.873* (1.105)	1.088 (1.008)
N	2160	7440	9600	2160	7440	9600
Pseudo R ²	0.285	0.321	0.322	0.142	0.203	0.189

Table 4 Determinants of punishment. Values in columns 1-3 reflect estimates from logistic models fitted to the decisions to punish (0: no deduction points assigned; 1: at least one deduction point assigned). Values in columns 3-6 reflect effect estimates from left-censored Tobit models fitted to the number of deduction points assigned. ‘Target’s contribution’ is the contribution of the punished participant; ‘Punisher’s contribution’ is the contribution of the participant punishing; ‘Average contribution others’ is the mean contribution of the other two members of the group; ‘Received punishment in *t*-1’ is the punishment amount received from others in the previous period; ‘Period’ is the period number; ‘Final period’ is a dummy for the last period; ‘MTurk’ is a dummy for the MTurk sample. Robust standard errors clustered on group. * $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$.

also reveal that both in the laboratory and online the positive effects of ‘received punishment in *t*-1’ on punishment reported in Table 4 are due to instances of anti-social punishment.

In sum, our results show that basic patterns of cooperation and punishment behavior in the laboratory are largely replicable online, and thus robust to changes in the experimental method. Participants in our online sample initially contribute more but, in the absence of punishment opportunities, cooperation decays at similar rates. In both samples, peer punishment is mainly pro-social in nature, and its introduction increases and stabilizes cooperation.

5 Attrition: endogenous or exogenous?

Our observation that, across conditions, experimental results are quite robust already suggests that endogenous attrition due to what has happened in the experiment (and hence selection) is not a big issue in our data. Here, we investigate this issue more rigorously.

Crucially, we find no evidence that attrition was *selective* in our experiment: dropout rates did not vary with the experimental condition (absence or presence of punishment). Table 5 details the results of a proportional hazards model (Jenkins 1995) fitted to instances where participants dropped out. The first model confirms the visual impression from Figure 1 that dropout rates are relatively high in the first period of each experimental condition (columns 1-3). Most likely, this effect is due to the fact that participants have to wait for their group to form (before the start of the first condition), or for all members of their group to finish reading the instructions and completing the comprehension questions (before the start of the second condition). Over the course of each of the conditions, the attrition rates slightly decrease (‘period’ has a negative estimate), suggesting that over time participants become more loyal to the task.

Models 3-5 show that attrition is much more likely when a group member has dropped out in the previous period. This seems indicative of ‘cascading inattention’: when a participant drops out of the session (e.g. due to inactivity, a closed connection, or waning attention), their group members will have to wait for some time before they can proceed.²⁸ Reduced attention may lead to additional attrition.

The results in columns 4 and 5 of Table 5 indicate that attrition did not depend on cooperation and punishment behavior in each of the experimental conditions. Specifically, dropouts did not depend on the behavior of the dropped-out participants relative to their fellow group mates, or on their earnings. The model in column 4 (fitted to the data from the experiment without punishment) shows that dropouts did not depend on relative average contributions. In addition to that, the model in column 5 (fitted to the data from the experiment with punishment) reveals that individuals who dropped out had neither received more punishment relative to their group members who did not drop out, nor did they differ from them in terms of the punishment towards others.

²⁸ In our software a subject was considered to have dropped out after 20 seconds of inactivity. This delay was introduced in order to allow people with brief connection irregularities to still complete the experiment.

<i>Participant's drop out in period t (0=no; 1=yes)</i>					
	<i>Pooled data</i>			<i>Without punishment</i>	<i>With punishment</i>
	(1)	(2)	(3)	(4)	(5)
Punishment available	0.056 (0.598)	0.362 (0.612)	0.107 (0.611)		
Period	-0.093* (0.051)	-0.118** (0.053)	-0.094* (0.053)	-0.265*** (0.080)	-0.150* (0.082)
First period	2.484*** (0.377)	2.375*** (0.376)	2.554*** (0.382)		
Earnings		-0.002 (0.143)	0.011 (0.143)		
Group member(s) dropped out in previous period			1.890*** (0.382)	3.636*** (0.394)	2.034*** (0.573)
Relative average contribution				0.010 (0.033)	-0.082 (0.053)
Relative average punishment received					-0.104 (0.092)
Relative average punishment given					-0.204 (0.214)
Constant	-4.064*** (0.317)	-3.979*** (0.318)	-4.220*** (0.328)	-3.519*** (0.466)	-4.282*** (0.472)
N	8334	8327	8332	3539	3527
AIC	893.56	877.20	860.27	325.98	302.30

Table 5 Determinants of attrition. Values reflect estimates from proportional hazards models fitted to binary events of participants staying (0) or dropping out (1) in a given round of the session, conditional on not having dropped out yet. ‘Punishment available’ is a dummy for the presence or absence of punishment; ‘Period’ is the period number; ‘First period’ is a dummy for the first period; ‘Earnings’ reflect participants’ total earnings relative to all other participants in the experiment in a given period; ‘Group member(s) dropped out in previous period’ is a dummy taking the value of 0 (1) when none (at least one) of the group members had left the session in the previous round (potentially delaying the progress within the session); ‘Relative average contribution’ is the participant’s average contribution to the public good minus the average contribution of their fellow group members in all rounds of the session so far; ‘Relative average received (given) punishment’ are the average punishment received (given) by a participant minus the average punishment received (given) by their fellow group members in all rounds so far. * $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$.

6 Discussion

In light of the results presented here, one might feel tempted to embrace interactive online experimentation as a valuable complement to laboratory studies – and others might even see it as a cost-efficient substitute. The measures presented here address the most important methodological issues for conducting interactive experiments online, and our case study illustrates that established results from the laboratory can be replicated online. However, future research needs to establish how generalizable our results are to other research questions as, despite these measures, methodological differences between laboratory and online experimentation inevitably remain.

For instance, depending on the nature of the experiment, online participants can conceivably communicate with each other to share their knowledge, strategies and even experimental materials more quickly than their laboratory counterparts. As mentioned earlier, most of the

forums that monitor the online community have mechanisms in place that prohibit the dissemination of materials, and participants themselves might find this practice prohibitively costly. Yet, one cannot completely rule out this possibility as laboratory and online participants can simply discuss an experiment through other channels. To some extent, the nature of interactive designs prevents participants from crafting intricate strategies beforehand, but this might not be the case for experimental designs where participants can figure out “correct” answers, and they might be at risk of being ineffective (Haigh 2016), or exhibiting reduced effect sizes (Chandler et al. 2015).

Comparisons between online and laboratory experiments can also be affected by differences in selection bias. Participants in online and laboratory experiments may self-select based on their opportunity cost of working time and their reservation wage. Indeed, opportunity costs and reservation wages might well differ between sessions conducted in the laboratory and online: laboratory participants might decide on whether or not to participate in an experiment by looking at the show-up fee paid and the travel costs they would incur (e.g. walk a long distance, experience bad weather, or even get dressed!), whereas for online participants such costs would typically be negligible. Interestingly though, results from Anderson et al. (2013) show that a comparable type of selection is unlikely to bias inference about the prevalence of other-regarding preferences. Thus, we have reasons to believe that our design is not particularly affected by the relatively low opportunity costs, but that others might be.

In this study we systematically controlled for what we think are the most daunting logistical issues for running an interactive experiment online. However, one could argue that some important methodological differences between laboratory and online experiments still remain, and that such discrepancies may potentially affect findings and treatment comparisons, regardless of experimental designs being interactive or not. For instance, we replicate classic patterns of behavior in an environment with less control but also find an important disparity between initial contributions. As we pointed out earlier though, this divergence is consistent with the one found in related studies comparing different adult populations with college students using various recruitment methods (Belot et al. 2010; Carpenter and Seki 2011; Gächter et al. 2004; Stoop et al. 2012; Gächter and Herrmann 2011). We acknowledge that unobserved methodological differences might account for some of the variability observed in our results, but not all. Assessing the extent and severity of such discrepancies in other designs is certainly a topic for future research.

The similarities between our online and laboratory results suggest that interactive designs conducted over the internet can be robust to changes in the experimental method. Yet, our results are based on a comparison between sessions that differed in terms of both the experimental method (online versus in the physical laboratory) and the subject pool (MTurk workers versus university participants). Further assessment of online experimentation could include a systematic study of the isolated effects of the method for collecting data (online or laboratory) on the one hand, and the subject pool (‘workers’ from an online labor market or university students) on the other hand, e.g. by running online experiments with university students or inviting MTurkers into the physical laboratory.

In summary, we see our paper as a guide for researchers to think about relevant issues before deciding whether the online or the physical laboratory is most appropriate for their research question. Some might conclude that the loss of control is too big a problem for their designs, whereas others are willing to bear that loss of control. In the end, the extent of any loss of control is an empirical question and we encourage researchers to add to our first piece of experimental evidence.

7 Conclusion

In this paper we presented a detailed conceptual and methodological discussion of conducting interactive experiments in the physical laboratory and online. We illustrated similarities and differences using a repeated public goods experiment without and with punishment. Our comparative results suggest that online data quality is adequate and reliable, making online interactive experimentation a potentially valuable complement to laboratory studies.

Most importantly, attrition, though a significant nuisance in online experiments, did not compromise the internal validity of our data because attrition was unrelated to what happened in our experiment. Future research will need to establish how generalizable this result is to other interactive decision problems, in particular when attrition might be treatment-specific, which poses the biggest problem to internal validity (Zhou and Fishbach 2016). Future research should also investigate how individual characteristics of participants (e.g. social preferences) and aspects of the experimental design (e.g. group size, number of periods, complexity of instructions) affect dropouts.

We observed that cooperation levels in our online sample are substantially higher than in the laboratory, and are on the high end of the range of cooperation levels observed in the cross-cultural samples of Herrmann et al. (2008). These differences can be partly (but not completely) explained by the age of MTurkers relative to students in typical laboratory samples. Still, it is unclear whether some other differences in terms of the participants' demographics, the perceived degree of anonymity, or the degree of familiarity with the experimental paradigm influence our results. We believe that future research should explore such avenues.

Acknowledgements We thank Thorsten Chmura, David Cooper, Dennie van Dolder, Urs Fischbacher, Marcus Giamattei, Roberto Hernan Gonzalez, Oliver Hauser, Susie Jackson, Daniele Nosenzo, Charles Noussair, David G. Rand, Ruslan Shichman, Katrin Schmelz, Jonathan Schulz, Richard Upward, Ori Weisel, Irenaeus Wolff, participants at the 2014 ESA Meeting in Ft. Lauderdale, and two anonymous referees for their helpful comments. We thank the European Research Council (ERC) Advanced Investigator Grant 295707 and the ESRC Network for Integrated Behavioural Science (NIBS, ES/K002201/1) for financial support.

References

- Abeler, J., & Nosenzo, D. (2015). Self-selection into laboratory experiments: pro-social motives versus monetary incentives. *Experimental Economics*, 18(2), 195-214, doi:10.1007/s10683-014-9397-9.
- Amir, O., Rand, D. G., & Gal, Y. K. (2012). Economic Games on the Internet: The Effect of \$1 Stakes. *Plos One*, 7(2), doi:ARTN e3146110.1371/journal.pone.0031461.
- Anderhub, V., Muller, R., & Schmidt, C. (2001). Design and evaluation of an economic experiment via the Internet. *Journal of Economic Behavior & Organization*, 46(2), 227-247, doi:Doi 10.1016/S0167-2681(01)00195-0.
- Anderson, J., Burks, S. V., Carpenter, J., Gotte, L., Maurer, K., Nosenzo, D., et al. (2013). Self-selection and variations in the laboratory measurement of other-regarding preferences across subject pools: evidence from one college student and two adult samples. *Experimental Economics*, 16(2), 170-189, doi:10.1007/s10683-012-9327-7.
- Arechar, A., Kraft-Todd, G., & Rand, D. G. (2015). Notes from a Day on the Forums: Recommendations for Maintaining a Good Reputation as an Amazon Mechanical Turk Requester. *Social Science Research Network*, 2602075.
- Behrend, T. S., Sharek, D. J., Meade, A. W., & Wiebe, E. N. (2011). The viability of crowdsourcing for survey research. *Behavior Research Methods*, 43(3), 800-813, doi:10.3758/s13428-011-0081-0.
- Belot, M., Duch, R., & Miller, L. (2010). Who should be called to the lab? A comprehensive comparison of students and non-students in classic experimental games. *University of Oxford, Nuffield College Discussion Papers*.((2010-001)).
- Berinsky, A. J., Huber, G. A., & Lenz, G. S. (2012). Evaluating Online Labor Markets for Experimental Research: Amazon.com's Mechanical Turk. *Political Analysis*, 20(3), 351-368, doi:10.1093/pan/mpr057.
- Bock, O., Baetge, I., & Nicklisch, A. (2014). hroot: Hamburg Registration and Organization Online Tool. *European Economic Review*, 71, 117-120, doi:10.1016/j.eurocorev.2014.07.003.
- Buhrmester, M., Kwang, T., & Gosling, S. D. (2011). Amazon's Mechanical Turk: A New Source of Inexpensive, Yet High-Quality, Data? *Perspectives on Psychological Science*, 6(1), 3-5, doi:10.1177/1745691610393980.
- Carpenter, J., & Seki, E. (2011). Do Social Preferences Increase Productivity? Field Experimental Evidence from Fishermen in Toyama Bay. *Economic Inquiry*, 49(2), 612-630, doi:10.1111/j.1465-7295.2009.00268.x.
- Chandler, J., Mueller, P., & Paolacci, G. (2014). Nonnaïveté among Amazon Mechanical Turk workers: Consequences and solutions for behavioral researchers. *Behavior Research Methods*, 46(1), 112-130, doi:10.3758/s13428-013-0365-7.
- Chandler, J., Paolacci, G., Peer, E., Mueller, P., & Ratliff, K. A. (2015). Using Nonnaive Participants Can Reduce Effect Sizes. *Psychological Science*, 26(7), 1131-1139.
- Chesney, T., Chuah, S. H., & Hoffmann, R. (2009). Virtual world experimentation: An exploratory study. *Journal of Economic Behavior & Organization*, 72(1), 618-635, doi:10.1016/j.jebo.2009.05.026.

- Cleave, B. L., Nikiforakis, N., & Slonim, R. (2013). Is there selection bias in laboratory experiments? The case of social and risk preferences. *Experimental Economics*, 16(3), 372-382, doi:10.1007/s10683-012-9342-8.
- Difallah, D. E., Catasta, M., Demartini, G., Ipeirotis, P. G., & Cudre-Mauroux, P. (2015). The Dynamics of Micro-Task Crowdsourcing The Case of Amazon MTurk. *Www'15 Companion: Proceedings of the 24th International Conference on World Wide Web*, 617-617, doi:10.1145/2740908.2744109.
- Egas, M., & Riedl, A. (2008). The economics of altruistic punishment and the maintenance of cooperation. *Proceedings of the Royal Society B-Biological Sciences*, 275(1637), 871-878, doi:10.1098/rspb.2007.1558.
- Fehr, E., & Gächter, S. (2000). Cooperation and punishment in public goods experiments. *American Economic Review*, 90(4), 980-994, doi:DOI 10.1257/aer.90.4.980.
- Fehr, E., & Gächter, S. (2002). Altruistic punishment in humans. *Nature*, 415(6868), 137-140.
- Fischbacher, U. (2007). z-Tree: Zurich toolbox for ready-made economic experiments. *Experimental Economics*, 10(2), 171-178, doi:10.1007/s10683-006-9159-4.
- Gächter, S., & Herrmann, B. (2011). The limits of self-governance when cooperators get punished: Experimental evidence from urban and rural Russia. *European Economic Review*, 55(2), 193-210, doi:10.1016/j.euroecorev.2010.04.003.
- Gächter, S., Herrmann, B., & Thoni, C. (2004). Trust, voluntary cooperation, and socio-economic background: survey and experimental evidence. *Journal of Economic Behavior & Organization*, 55(4), 505-531, doi:10.1016/j.jebo.2003.11.006.
- Gallo, E., & Yan, C. (2015). The effects of reputational and social knowledge on cooperation. *Proceedings of the National Academy of Sciences of the United States of America*, 112(12), 3647-3652, doi:10.1073/pnas.1415883112.
- Goodman, J. K., Cryder, C. E., & Cheema, A. (2013). Data Collection in a Flat World: The Strengths and Weaknesses of Mechanical Turk Samples. *Journal of Behavioral Decision Making*, 26(3), 213-224, doi:10.1002/bdm.1753.
- Greiner, B. (2015). Subject pool recruitment procedures: organizing experiments with ORSEE. [journal article]. *Journal of the Economic Science Association*, 1(1), 114-125, doi:10.1007/s40881-015-0004-4.
- Guillen, P., & Veszteg, R. F. (2012). On "lab rats". *Journal of Socio-Economics*, 41(5), 714-720, doi:10.1016/j.socec.2012.07.002.
- Haigh, M. (2016). Has the Standard Cognitive Reflection Test Become a Victim of Its Own Success? *Advances in Cognitive Psychology*, 12(3), 145-149, doi:10.5709/acp-0193-5.
- Hauser, O. P., Hendriks, A., Rand, D. G., & Nowak, M. A. (2016). Think global, act local: Preserving the global commons. [Article]. *Scientific reports*, 6, 36079, doi:10.1038/srep36079.
- Hergueux, J., & Jacquemet, N. (2015). Social preferences in the online laboratory: a randomized experiment. *Experimental Economics*, 18(2), 251-283, doi:10.1007/s10683-014-9400-5.
- Herrmann, B., Thoni, C., & Gächter, S. (2008). Antisocial punishment across societies. *Science*, 319(5868), 1362-1367.
- Horton, J. J., Rand, D. G., & Zeckhauser, R. J. (2011). The online laboratory: conducting experiments in a real labor market. *Experimental Economics*, 14(3), 399-425, doi:10.1007/s10683-011-9273-9.
- Ipeirotis, P. G. (2010). Analyzing the Amazon Mechanical Turk Marketplace. *XRDS: Crossroads, The ACM Magazine for Students*, 17(2), 16-21.

- Jenkins, S. P. (1995). Easy Estimation Methods for Discrete-Time Duration Models. *Oxford Bulletin of Economics and Statistics*, 57(1), 129-138, doi:DOI 10.1111/j.1468-0084.1995.tb00031.x.
- Klein, R. A., Ratliff, K. A., Vianello, M., Adams, R. B., Bahnik, S., Bernstein, M. J., et al. (2014). Investigating Variation in Replicability A "Many Labs" Replication Project. *Social Psychology*, 45(3), 142-152, doi:10.1027/1864-9335/a000178.
- Krawczyk, M. (2011). What brings your subjects to the lab? A field experiment. *Experimental Economics*, 14(4), 482-489, doi:10.1007/s10683-011-9277-5.
- Krupnikov, Y., & Levine, A. S. (2014). Cross-sample comparisons and external validity. *Journal of Experimental Political Science*, 1(1), 59.
- List, J. A. (2004). Young, selfish and male: Field evidence of social preferences. *Economic Journal*, 114(492), 121-149, doi:DOI 10.1046/j.0013-0133.2003.00180.x.
- Mullinix, K. J., Leeper, T. J., Druckman, J. N., & Freese, J. (2015). The Generalizability of Survey Experiments. *Journal of Experimental Political Science*, 2(02), 109-138, doi:doi:10.1017/XPS.2015.19.
- Nishi, A., Shirado, H., Rand, D. G., & Christakis, N. A. (2015). Inequality and visibility of wealth in experimental social networks. *Nature*, 526(7573), 426+, doi:10.1038/nature15392.
- Paolacci, G., & Chandler, J. (2014). Inside the Turk: Understanding Mechanical Turk as a Participant Pool. *Current Directions in Psychological Science*, 23(3), 184-188, doi:10.1177/0963721414531598.
- Paolacci, G., Chandler, J., & Ipeirotis, P. G. (2010). Running experiments on Amazon Mechanical Turk. *Judgment and Decision Making*, 5(5), 411-419.
- Peer, E., Samat, S., Brandimarte, L., & Acquisti, A. (2016). Beyond the Turk: An Empirical Comparison of Alternative Platforms for Crowdsourcing Online Behavioral Research. *Social Science Research Network*, 2594183.
- Peer, E., Vosgerau, J., & Acquisti, A. (2014). Reputation as a sufficient condition for data quality on Amazon Mechanical Turk. *Behavior Research Methods*, 46(4), 1023-1031, doi:10.3758/s13428-013-0434-y.
- Pfattheicher, S., & Schindler, S. (2015). Understanding the Dark Side of Costly Punishment: The Impact of Individual Differences in Everyday Sadism and Existential Threat. *European Journal of Personality*, 29(4), 498-505, doi:10.1002/per.2003.
- Rand, D. G. (2012). The promise of Mechanical Turk: How online labor markets can help theorists run behavioral experiments. *Journal of Theoretical Biology*, 299, 172-179, doi:10.1016/j.jtbi.2011.03.004.
- Rand, D. G., Peysakhovich, A., Kraft-Todd, G. T., Newman, G. E., Wurzbacher, O., Nowak, M. A., et al. (2014). Social heuristics shape intuitive cooperation. *Nature Communications*, 5, doi:ARTN 367710.1038/ncomms4677.
- Reips, U.-D. (2000). The Web experiment method: Advantages, disadvantages, and solutions. In M. H. Birnbaum (Ed.), *Psychological experiments on the Internet* (pp. 89-117). San Diego, CA: Academic Press.
- Schmelz, K., & Ziegelmeyer, A. (2015). Social distance and control aversion: Evidence from the Internet and the laboratory. *Research Paper Series Thurgau Institute of Economics and Department of Economics at the University of Konstanz, TWI-RPS 100*.

- Stagnaro, N., Arechar, A., & Rand, D. G. (2016). From Good Institutions to Good Norms: Top-Down Incentives to Cooperate Foster Prosociality but not Norm Enforcement. *Social Science Research Network*, 2720585.
- Stewart, N., Ungemach, C., Harris, A. J. L., Bartels, D. M., Newell, B. R., Paolacci, G., et al. (2015). The average laboratory samples a population of 7,300 Amazon Mechanical Turk workers. *Judgment and Decision Making*, 10(5), 479-491.
- Stoop, J., Noussair, C. N., & Van Soest, D. (2012). From the Lab to the Field: Cooperation among Fishermen. *Journal of Political Economy*, 120(6), 1027-1056, doi:10.1086/669253.
- Suri, S., & Watts, D. J. (2011). Cooperation and Contagion in Web-Based, Networked Public Goods Experiments. *Plos One*, 6(3), doi:ARTN e1683610.1371/journal.pone.0016836.
- Wang, J., Suri, S., & Watts, D. J. (2012). Cooperation and assortativity with dynamic partner updating. *Proceedings of the National Academy of Sciences of the United States of America*, 109(36), 14363-14368, doi:10.1073/pnas.1120867109.
- Weinberg, J. D., Freese, J., & McElhattan, D. (2014). Comparing data characteristics and results of an online factorial survey between a population-based and a crowdsourced-recruited sample. *Sociological Science*, 1, 292-310.
- Zhou, H. T., & Fishbach, A. (2016). The Pitfall of Experimenting on the Web: How Unattended Selective Attrition Leads to Surprising (Yet False) Research Conclusions. *Journal of personality and social psychology*, 111(4), 493-504, doi:10.1037/pspa0000056.

Appendix A – Additional statistical analyses

<i>Contribution to the public good</i>								
	Tobit estimation				Multilevel mixed effects estimation			
	No punishment (left censored)		Punishment (right censored)		No punishment		Punishment	
Period	-1.037*** (0.160)	-1.042*** (0.159)	0.682** (0.282)	0.711** (0.281)	-0.731*** (0.043)	-0.737*** (0.044)	0.091*** (0.024)	0.093*** (0.025)
Final period	-2.512*** (0.881)	-2.611*** (0.907)	-5.795*** (1.797)	-5.864*** (1.863)	-1.277*** (0.408)	-1.315*** (0.420)	-1.119*** (0.233)	-1.082*** (0.241)
MTurk	5.421*** (1.867)	5.390** (2.130)	4.193 (4.904)	8.282 (5.919)	4.218*** (1.394)	5.084*** (1.466)	0.893 (1.249)	2.009 (1.320)
Age		0.120** (0.058)		-0.003 (0.144)		0.028 (0.029)		-0.029 (0.025)
Female		2.370** (0.976)		-4.072 (2.595)		0.927* (0.516)		0.118 (0.438)
Single child		0.730 (1.262)		-1.096 (4.071)		0.051 (0.715)		-0.245 (0.607)
Foreign		-1.042 (2.557)		6.286 (4.622)		0.801 (1.196)		0.401 (1.065)
Membership		1.242 (1.452)		5.536 (4.661)		1.764*** (0.651)		1.520*** (0.556)
Constant	11.402*** (1.650)	6.544** (2.778)	29.601*** (4.232)	26.270*** (6.756)	12.451*** (1.246)	10.043*** (1.530)	15.871*** (1.109)	15.368*** (1.348)
N	3200	3050	3200	3050	3200	3050	3200	3050
Chi ² / F	34.45	15.35	3.75	1.96	503.80	496.65	25.67	32.03

Table A1 Cooperation dynamics. Tobit and Multilevel mixed effects estimation, which allows for individual and group differences, as well as for treatment-specific residuals. ‘Period’ is period number; ‘Final period’ is a dummy for last period; ‘MTurk’ is a dummy for the MTurk sample; ‘Age’ is the participant’s age; ‘Female’ is a dummy for female participants; ‘Foreign’ is a dummy for participants who grew up outside the US; ‘Membership’ is a dummy for participants who were members of a social club. Robust standard errors clustered on groups for the Tobit model; * $p < 0.1$, ** $p < 0.05$, *** $p < 0.01$.

<i>Contribution to the public good (no punishment)</i>				
	Tobit estimation		Multilevel mixed effects estimation	
	Period	-0.485*** (0.085)	-0.492*** (0.090)	-0.500*** (0.052)
Final period	-1.600** (0.757)	-1.741** (0.781)	-0.812** (0.394)	-0.864** (0.407)
Mean peer contribution in $t-1$	1.027*** (0.054)	1.000*** (0.058)	0.461*** (0.027)	0.460*** (0.027)
MTurk	0.759 (0.778)	1.827 (1.142)	2.252*** (0.862)	2.980*** (0.985)
Age		0.033 (0.039)		0.036 (0.031)
Female		2.022*** (0.674)		1.230** (0.557)
Single child		0.128 (1.017)		0.128 (0.789)
Foreign		-0.321 (1.903)		0.191 (1.315)
Membership		2.153** (0.975)		1.692** (0.690)
Constant	-0.696 (0.931)	-3.582** (1.599)	7.073*** (0.870)	4.412*** (1.296)
N	2880	2745	2880	2745
Chi ² / F	177.16	80.17	822.66	799.84

Table A2 Cooperation dynamics (no punishment). Left-censored Tobit and multilevel mixed effects estimation, which allows for individual and group differences, as well as for treatment-specific residuals. ‘Period’ is period number; ‘Final period’ is a dummy for the last period; ‘Mean peer contribution in $t-1$ ’ is the average contribution of the other members in the group in $t-1$; ‘MTurk’ is a dummy for the MTurk sample; demographic controls are the same of Table A1. Robust standard errors clustered on groups; Robust standard errors clustered on groups for the Tobit model; * $p < 0.1$, ** $p < 0.05$, *** $p < 0.01$.

Decision to punish (0=no; 1=yes)

	Logit estimation								Multilevel mixed effects estimation							
	<i>Prosocial punishment</i>				<i>Antisocial punishment</i>				<i>Prosocial punishment</i>				<i>Antisocial punishment</i>			
	Lab	MTurk	Pooled	Pooled & controls	Lab	MTurk	Pooled	Pooled & controls	Lab	MTurk	Pooled	Pooled & controls	Lab	MTurk	Pooled	Pooled & controls
Punisher's contribution	0.046 (0.051)	-0.003 (0.023)	0.011 (0.022)	-0.002 (0.026)	-0.143*** (0.024)	-0.153*** (0.026)	-0.142*** (0.018)	-0.130*** (0.020)	0.000 (0.006)	0.003 (0.005)	0.002 (0.004)	0.000 (0.004)	-0.008*** (0.001)	-0.004*** (0.001)	-0.005*** (0.000)	-0.004*** (0.000)
Target's contribution	-0.059 (0.038)	-0.118*** (0.029)	-0.100*** (0.023)	-0.097*** (0.023)	-0.013 (0.025)	0.010 (0.021)	0.002 (0.015)	0.004 (0.020)	-0.012*** (0.004)	-0.021*** (0.003)	-0.018*** (0.002)	-0.018*** (0.002)	0.000 (0.002)	0.000 (0.001)	0.000 (0.001)	0.000 (0.001)
Others' avg. contribution	0.038 (0.035)	0.054*** (0.019)	0.047*** (0.017)	0.065*** (0.017)	-0.025 (0.028)	-0.019 (0.024)	-0.023 (0.018)	-0.017 (0.019)	0.013*** (0.005)	0.011*** (0.003)	0.012*** (0.003)	0.015*** (0.003)	0.000 (0.001)	-0.001 (0.001)	-0.001 (0.001)	0.000 (0.001)
Received <i>p</i> in t-1	-0.034 (0.027)	-0.025 (0.045)	-0.028 (0.023)	-0.076** (0.032)	0.129*** (0.031)	0.069** (0.034)	0.094*** (0.023)	0.103*** (0.025)	-0.009 (0.005)	-0.004 (0.006)	-0.007 (0.004)	-0.012*** (0.004)	0.010*** (0.001)	0.003*** (0.001)	0.004*** (0.001)	0.004*** (0.001)
Period	-0.117*** (0.039)	-0.111** (0.044)	-0.113*** (0.032)	-0.103*** (0.032)	0.024 (0.062)	0.036 (0.056)	0.031 (0.041)	0.032 (0.048)	-0.014 (0.009)	-0.011** (0.005)	-0.012*** (0.004)	-0.013*** (0.004)	0.004* (0.002)	0.001 (0.001)	0.001 (0.001)	0.001 (0.001)
Final period	-0.115 (0.341)	-0.077 (0.366)	-0.077 (0.275)	-0.236 (0.313)	-1.800*** (0.587)	0.595* (0.340)	-0.159 (0.388)	-0.014 (0.434)	-0.037 (0.086)	-0.045 (0.047)	-0.043 (0.041)	-0.040 (0.042)	-0.081*** (0.019)	0.018*** (0.006)	0.010* (0.006)	0.011* (0.006)
MTurk			-0.855*** (0.307)	-0.983*** (0.302)			-1.132*** (0.347)	-1.983*** (0.403)			-0.190*** (0.059)	-0.215*** (0.072)		-0.044*** (0.016)	-0.088*** (0.018)	
Age				0.005 (0.017)			0.025 (0.019)					0.001 (0.003)			0.001 (0.001)	
Female				-1.171*** (0.260)			-0.078 (0.381)					-0.204*** (0.052)			0.001 (0.011)	
Single child				-0.320 (0.376)			0.099 (0.457)					-0.085 (0.076)			0.002 (0.016)	
Foreign				1.403* (0.739)			-2.380*** (0.628)					0.154 (0.105)			-0.071** (0.028)	
Membership				0.019 (0.288)			-0.838* (0.462)					-0.070 (0.063)			-0.035** (0.014)	
Constant	0.206 (0.812)	0.202 (0.393)	0.854* (0.466)	1.216* (0.712)	-0.560 (0.379)	-2.137*** (0.682)	-0.878** (0.442)	-1.101 (0.682)	0.611*** (0.123)	0.437*** (0.092)	0.622*** (0.084)	0.700*** (0.122)	0.187*** (0.042)	0.104*** (0.016)	0.156*** (0.020)	0.176*** (0.027)
N	370	900	1270	1201	1790	6540	8330	7949	370	900	1270	1201	1790	6540	8330	7949
Chi ²	27.51	38.98	48.60	109.88	111.55	85.36	187.39	254.26	24.64	98.78	129.04	168.92	138.89	116.77	197.44	226.82

Table A3 Determinants of peer punishment. Logit and Multilevel mixed effects estimation, which allows for individual and group differences, as well as for treatment-specific residuals. We split the analysis into two different types of punishment. Pro-social punishment includes instances where the punisher's contribution to the public good in that round exceeded that of their target. Anti-social punishment includes instances where the target contributed at least as much as the punisher. 'Punisher's contribution' is the contribution of the participant punishing; 'Target's contribution' is the contribution of the punished participant; 'Average contribution others' is the mean contribution of the other two members of the group; 'Received punishment in t-1' is the punishment amount received from others in the previous period; 'Period' is the period number; 'Final period' is a dummy for the last period; 'MTurk' is a dummy for the MTurk sample; demographic controls are the same of Table A1. Robust standard errors clustered on groups for the Logit model * $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$.

Online Appendix

Appendix B – Experimental Instructions

Below is the description of the HIT of our MTurk experiment, as seen by participants.

The screenshot shows the Amazon Mechanical Turk HIT preview interface. At the top, a blue header bar contains the text: "Make decisions and earn \$1 for answering a quiz correctly plus up to \$9.80 in bonus." Below this, a white bar displays the following information: "Requester: Studies on MTurk", "Reward: \$1.00 per HIT", "HITs available: 0", and "Duration: 45 Minutes". A light blue bar below that states "Qualifications Required: HIT Approval Rate (%) for all Requesters' HITs greater than or equal to 90, Location is US". The main content area is titled "HIT Preview" and contains the following text:

Title of the study: Group decision making.
Length of the study: about 30 minutes.
Payment: A flat fee of \$1.00 upon completion of a quiz plus a bonus that depends on your decisions and the decisions of others.
What this study is about
This is a study about group decision making. You will make decisions for which you will receive a bonus dependent on your and others' decisions, in addition of a flat payment of \$1 for completing a quiz.
Voluntary participation and confidentiality
Participation in this study is entirely voluntary, and may be withdrawn at any given moment without further consequences. All decisions will be made anonymously, and results will only be analyzed at the group level for publications in scientific journals.
Browser compatibility
This HIT will not work on Internet Explorer.
Risks
To our knowledge, there are no risks involved in participating in this survey.
Informed consent
By accepting this HIT you give us informed consent that we can use your answers in anonymized form for research purposes only.

- Groups of 4 real people recruited via MTurk
- 30 MINUTES REQUIRED **WITHOUT INTERRUPTION**, STARTING IMMEDIATELY.
- We can only allow up to 20 participants in this group of HITs
- We recommend that you click on the link below before accepting the HIT
- At the end of the experiment, you will receive a code to paste into the box below to receive credit for taking our experiment.

Survey link: [Click here](#)

Provide the survey code here:

Moreover, below are the instructions given to all of our participants (laboratory and MTurk). Each box corresponds to the screens shown.

Instructions

In this experiment you will play a game with the same three people for 10 rounds. In these rounds, you can earn Points. At the end of the experiment these Points will be converted into real money (**50 Points = 1 Dollar**).

Your task

At the beginning of each round, each participant receives 20 Points. You have to decide how many of the 20 Points you want to **contribute to a group project**. The other three members of your group make this decision at the same time.

The Points you do not contribute, you **keep for yourself**. These Points are added to your total. After all group members have made their decision, all Points contributed to the **group project** are added up, and this number of Points is **multiplied by 1.6**. The resulting number of Points is then divided equally among the group members (irrespective of how much they individually contributed to the group project).

In summary:
Your income in a round =
The Points you keep for yourself
plus
The Points you receive from the group project

Group project - Example 1

- All four players contribute 20 Points to the group project.
- Sum of contributions is 80 Points.
- This amount is multiplied by 1.6, resulting in 128 Points.
- Each participant receives $(128/4=)$ 32 Points from the group project.
- Therefore, the income of each player is 32 Points.

Group project - Example 2

- Participants A, B and C contribute each 20 Points to the group project.
- Participant D contributes 0 Points.
- Sum of contributions is 60 Points.
- This amount is multiplied by 1.6, resulting in 96 Points.
- Each participant receives $(96/4=)$ 24 Points from the group project.
- Therefore, the income of Participants A, B and C is 24 Points.
- The income of Participant D is 44 Points (20 kept for himself *plus* 24 from the group project).

Please make your decision within the **time limit** shown on your screen.

After all members of your group have made their decision, the results of the round will be shown to you. Once all players in your group are finished, a new round will begin. Again, you will receive 20 Points to start with. After 10 rounds, this first experiment is over. A second and final experiment will follow. Further information about this second experiment will be provided later.

Please click the link below if you understood the instructions.

Before the experiment itself starts, a brief quiz will check whether you understand your task.

Quiz

Please answer all control questions. They serve as a test for your understanding of the experiment.

1. At the start of a round, each group member receives 20 Points. Suppose nobody (including you) contributes any Points to the group project. What is:

- a) Your income?
- b) The income of each of the other group members?

2. At the start of a round, each group member receives 20 Points. Suppose you contribute 20 Points to the group project. All other group members each contribute 20 Points to the group project. What is:

- a) Your income?
- b) The income of each of the other group members?

3. At the start of a round, each group member receives 20 Points. Suppose the other three group members contribute a total of 40 Points to the group project.

- a) What is your income if you contribute 0 Points to the group project?
- b) What is your income if you contribute 10 Points to the project?

The first experiment of this session is now over. A second and final experiment will follow.

Instructions for the second experiment

In this second experiment you will again play a game in the **same group** for **10 rounds**. In each round your task is **very similar** to your task in the first experiment, but we introduce **one change**. After all group members have made their contributions to the group project, you can **assign Deduction Points** to each of the other members of your group.

For each Deduction Point you assign, **3 Points** will be deducted from the total of the recipient, and **1 Point** will be deducted from your total. In each round, you can assign **between 0 and 10 Deduction Points** to each of the other members of your group.

Deduction Points - Example

- All members of your group have made their contribution to the group project.
 - As before, your screen displays the results of this round so far.
 - You assign the following Deduction Points to Participants A, B and C: 2, 0 and 3.
 - This reduces your income in this round by $(2+0+3=)$ 5 Points.
 - The other Participants assign a total of 4 Deduction Points to you.
 - This reduces your income in this round by $(4 \times 3=)$ 12 Points.
- In this round, your total income reduction due to Deduction Points is $(5+12=)$ 17 Points.

Please be aware that when many Deduction Points are assigned to you, your income in a round may become **negative**. In that case, your income in that round will be set to 0 (zero). However, you will **always** incur the costs of assigning Deduction Points to others.

After all members of your group have assigned their Deduction Points, the results of the round will be shown to you. Once all players in your group are finished, a new round will begin. After 10 rounds of this game, a brief questionnaire will follow concluding this session. Your payment for this session will consist of:

Your guaranteed \$10 participation fee
+ Your income from the 10 rounds of the first experiment
+ Your income from the 10 rounds of the upcoming second experiment

At the start of this second experiment, you will receive an **extra 25 Points** to compensate for potential losses due to Deduction Points.

Please click the link below if you understood the instructions. A brief quiz will check whether you understand your task before the experiment itself starts.

[I have read and understood the instructions. Continue!](#)

Quiz

Please answer all control questions. They serve as a test for your understanding of the experiment.

1. All members of your group have made their contributions to the group project. Suppose that you assign the following Deduction Points to the other members of your group: 2, 3 and 0. By many Points will your income be reduced by assigning Deduction Points in this round?

2. By many Points will your income be reduced if you assign 0 Deduction Points in total?

3. By how many Points will your income be reduced if you receive a total of 7 Deduction Points from the other members of your group?

[Submit](#)

[Go back to instructions of the upcoming experiment.](#)

Appendix C – Cooperation and punishment on MTurk from a wider perspective

Although aggregate dynamics of cooperation were similar across our samples, cooperation levels were appreciably higher online than in the laboratory, especially in the absence of punishment (Figure 2 of the main text). Indeed, in comparison to the wide spectrum of cooperation levels observed in a large cross-cultural sample (Herrmann et al. 2008), our online sample ranks at the very top (Figures C1 and C2). Specifically, in the absence of punishment our online sample tops the ranks (Fig. C1). In the presence of punishment it starts with the highest cooperation levels and ranks third in average cooperation (Fig. C2), while punishment is used less frequently than in any other sample (Fig. C3).

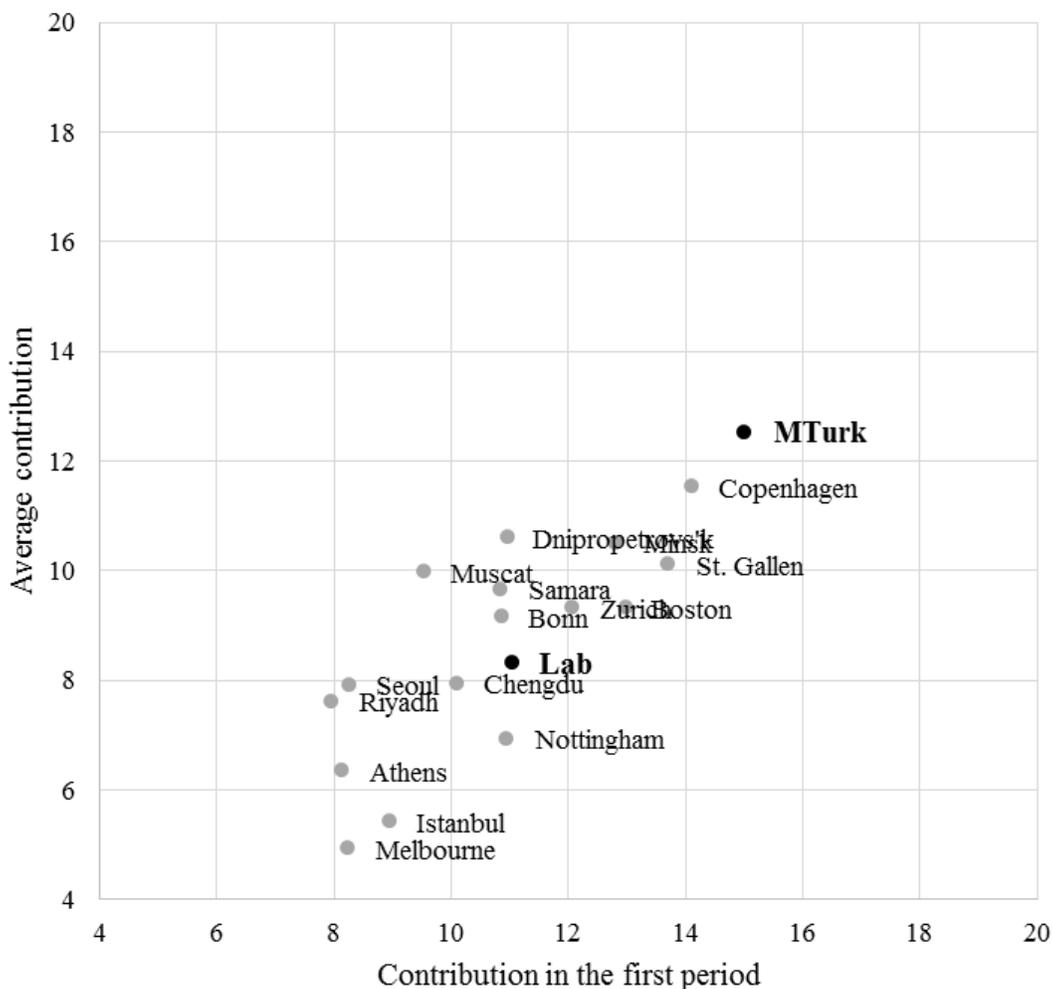


Fig. C1 Contributions in the first period and average contribution in the No-Punishment condition in our samples from MTurk, the laboratory, and the 16 cities included in Herrmann et al. (2008)

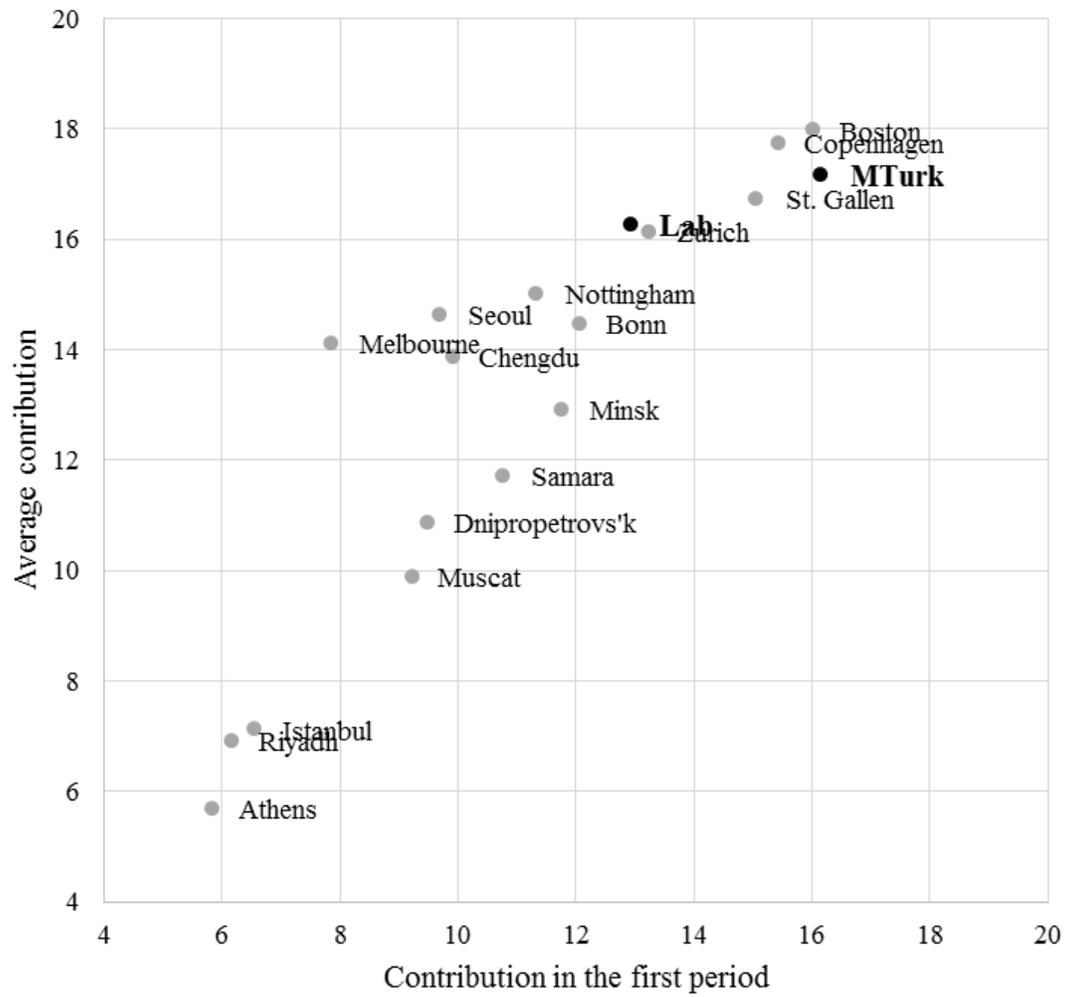


Fig. C2 Contributions in the first period and average contribution in the Punishment condition in our samples from MTurk, the laboratory, and the 16 cities included in Herrmann et al. (2008)

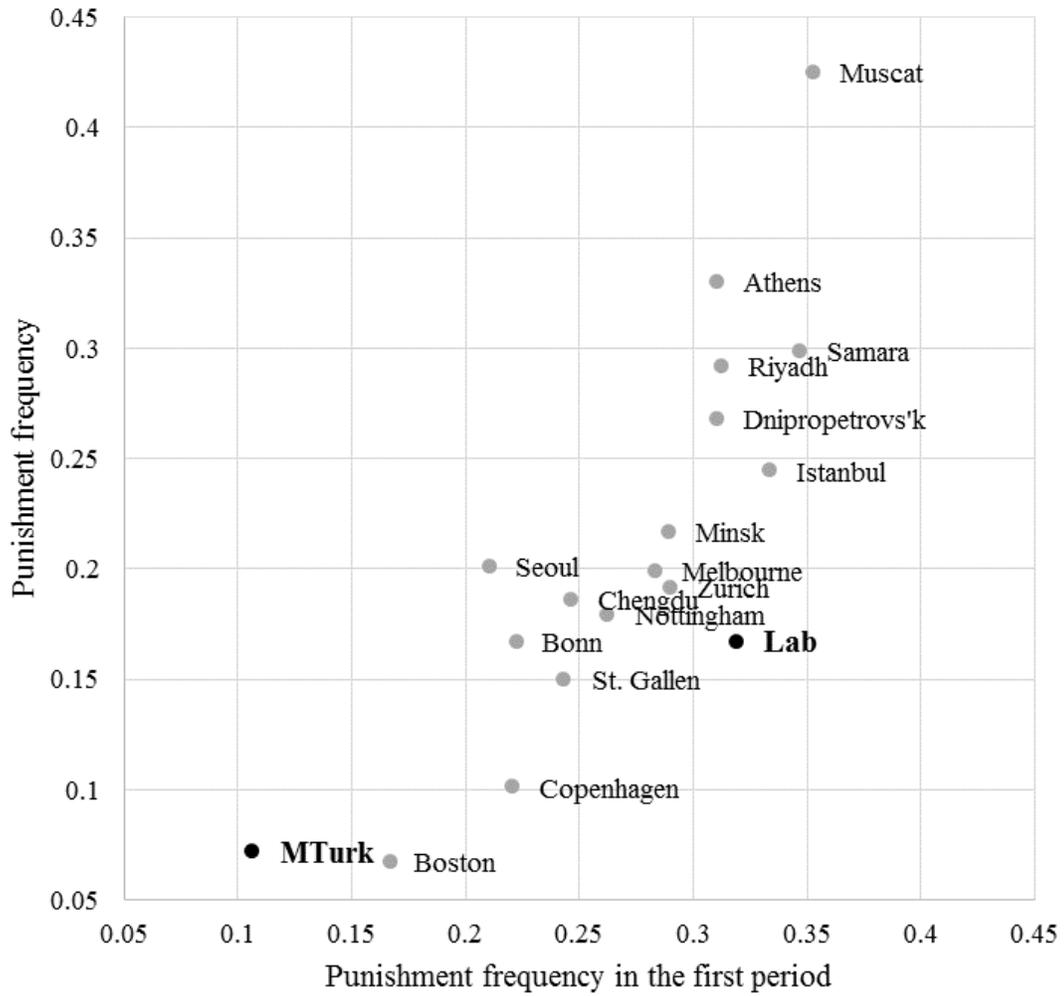


Fig. C3 Punishment frequencies in the first period and overall frequencies in our samples from MTurk, the laboratory, and the 16 cities included in Herrmann et al. (2008)