

Backhaus, Teresa; Breitmoser, Yves

Working Paper

God Does Not Play Dice, but Do We? On the Determinism of Choice in Long-run Interactions

Discussion Paper, No. 96

Provided in Cooperation with:

University of Munich (LMU) and Humboldt University Berlin, Collaborative Research Center Transregio 190: Rationality and Competition

Suggested Citation: Backhaus, Teresa; Breitmoser, Yves (2018) : God Does Not Play Dice, but Do We? On the Determinism of Choice in Long-run Interactions, Discussion Paper, No. 96, Ludwig-Maximilians-Universität München und Humboldt-Universität zu Berlin, Collaborative Research Center Transregio 190 - Rationality and Competition, München und Berlin

This Version is available at:

<https://hdl.handle.net/10419/185766>

Standard-Nutzungsbedingungen:

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

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

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

Terms of use:

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

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

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

God Does Not Play Dice, but Do We?

On the Determinism of Choice in Long-Run Interactions

Teresa Backhaus (WZB)
Yves Breitmoser (HU Berlin)

Discussion Paper No. 96

May 11, 2018

God does not play dice, but do we?

On the determinism of choice in long-run interactions

Teresa Backhaus*
BDPEMS and WZB

Yves Breitmoser
Humboldt University Berlin

May 8, 2018

Abstract

When do we cooperate and why? This question concerns one of the most persistent divides between “theory and practice”, between predictions from game theory and results from experimental studies. For about 15 years, theoretical analyses predict completely-mixed “behavior” strategies, i.e. strategic randomization rendering “when” and “why” questions largely moot, while experimental analyses seem to consistently identify pure strategies, suggesting long-run interactions are deterministic. Reanalyzing 145,000 decisions from infinitely repeated prisoner’s dilemma experiments, and using data-mining techniques giving pure strategies the best possible chance, we conclude that subjects play semi-grim behavior strategies similar to those predicted by theory.

JEL-Code: C72, C73, C92, D12

Keywords: Repeated game, Behavior, Tit-for-tat, Mixed strategy, Memory, Belief-free equilibrium, Laboratory experiment

*We thank Andrea Ariu, Kai Barron, Friedel Bolle, Steffen Huck, Johannes Leutgeb, Vlada Pleshcheva, Sebastian Schweighofer-Kodritsch, Roland Strausz, Robert Stüber, Georg Weizsäcker, Roel van Veldhuizen, and Joachim Winter, as well as audiences in Schwanenwerder and Berlin for many helpful comments. Financial support by the Deutsche Forschungsgemeinschaft (BR 4648/1 and CRC TRR 190) is gratefully acknowledged. Address Backhaus: Wissenschaftszentrum Berlin für Sozialforschung, Reichpietschufer 50, 10785 Berlin, Germany, email: teresa.backhaus@wzb.eu, phone: +49 30 25491 411. Address Breitmoser: Spandauer Str. 1, 10099 Berlin, Germany, email: yves.breitmoser@hu-berlin.de, phone: +49 30 2093 99408.

1 Introduction

One of the most dynamic research fields over the last two decades has been behavioral game theory: the econometric and theoretical analysis of laboratory games to align observed behavior with game-theoretical concepts. How should we think of beliefs, utilities and choice of subjects, and is it possible to explain their decisions as responses to incentives? In some classes of games, most notably auctions, behavior seems to be reasonably consistent with theory after accounting simply for risk aversion (Bajari and Hortacsu, 2005) or biased beliefs (Eyster and Rabin, 2005). In generic normal-form games involving dominated strategies behavior is captured after relaxing rational expectations (Costa-Gomes et al., 2001); in games without dominated strategies behavior tends to reflect mainly logistic errors in choice (Weizsäcker, 2003; Brunner et al., 2011); and in games involving the distribution of monetary benefits, interdependence of preferences seems to organize behavior (Fehr and Schmidt, 1999; Charness and Rabin, 2002). Particular behavioral models tend to be disputed, but overall, there has been substantial progress in aligning observed behavior and theoretical predictions across many classes of games.

A class of games that arguably experienced less progress in aligning behavior and predictions is the large class of repeated games. Repeated games are the main tool in modeling long-run interactions, in particular to study cooperation and defection, and have been a core object of game-theoretic analyses at least since the Folk Theorem for repeated games with discounting (Fudenberg and Maskin, 1986). Part of an explanation for the slow progress may be that theoretical predictions for repeated games tend to be less specific than for other games. The main reason, however, seems to be independent of that: Behavioral analyses of individual strategies in repeated games typically do not assume subgame perfection, even though this normally represents an indisputable assumption in theoretical and behavioral analyses of dynamic games (with the exception of Fudenberg and Levine, 1993). Relaxing subgame perfection is necessary to justify the widespread intuition that human behavior is directly reciprocal, as captured by tit-for-tat in repeated games. When taking tit-for-tat reciprocity as given, dropping subgame perfection in behavioral analyses of repeated games seems inevitable, which perhaps instills the belief that theoretical results maintaining subgame perfection cannot inform such analyses.

This state of affairs raises a provocative question: Who is right? *Theoretical analyses* studying subgame perfect equilibria, assuming subjects understand non-credible threats and seek robustness to, for example, imperfect monitoring (following Kandori, 2002, and Ely and Välimäki, 2002), or *behavioral analyses* analyzing rather complex deterministic choice rules such as tit-for-tat, “lenient grim”, or “tit-for-2-tats” (Fudenberg et al., 2012) that mostly violate subgame perfection but seek to build on psychological interpretations of behavior?

With this paper, we seek to tackle this question. Our basis for identification is arguably the main prediction of theoretical analyses of repeated games from the last

twenty years: Following Kandori (2002) and Ely and Välimäki (2002), a large body of theoretical work analyzed robustness to imperfect (private) monitoring and concluded that players seeking robustness are likely to play *behavior strategies* (see also Ely et al., 2005, and Mailath and Samuelson, 2006). Private monitoring is imperfect if opponents' actions are not observable but players get imperfect private signals about their opponents' actions. This covers the case that players fail to perfectly observe or remember their opponents' actions, which seems both empirically and behaviorally plausible. The main prediction of this literature is that subjects should play belief-free equilibria (or weakly belief-free equilibria, Kandori, 2011), i.e. they should randomize each round to ensure the indifference of opponents at all points of time—which implies that opponents do not need to be able to perfectly remember previous actions. This prediction diametrically opposes the assumption of deterministic choice rules made in most recent behavioral analyses (the exception, Breitmoser, 2015, is discussed shortly) and thus provides a powerful foundation for answering the question(s) raised above. The difficulty is that an enormous data set and specific econometric tools are required to conclusively discriminate between behavior strategies, where subjects strategically randomize, and pure strategies, where deviations simply are stochastic mistakes—substantially more data and different tools than in the suggestive analysis of Breitmoser (2015).

The data concern has been solved very recently by Dal Bó and Fréchette (2018), who compiled a very large data set comprising 12 previous experiments (with 32 treatments and around 150,000 decisions), but in their analysis they maintain the assumption of deterministic choice rules. This data set finally allows us to tackle the questions raised above. To conclusively do so, we use a technique that belongs to the realm of data mining: We select models from wide ranges of one- and two-memory pure strategies, as well as from a range of different “switching” rules to capture inconsistency at the individual level. These ranges comprise up to 10^{44} models and represent exhaustive mining for the best combination of pure strategies. Only the best such model will be evaluated against the simple three-parametric “semi-grim” behavior strategy¹ previously identified by Breitmoser (2015)—importantly without penalizing the best model for the selection steps that preceded the evaluation stage. This approach, mining for the best-possible model and then evaluating the model as if it had been hypothesized ex-ante, is called data dredging and is frowned upon for obvious reasons. It drastically biases p -values in favor of the mined model, but in our case this approach helps us to ensure that we are giving pure strategies the best possible chance. Still, we find that not even this exhaustively data-mined mixture improves on the simple hypothesis that subjects play behavior strategies with a rather specific structure (semi-grim)—consistently across experiments. Mostly, the mined model actually fits significantly worse.

We conclude that individual behavior is indeed best described by the memory-1

¹Subjects randomize on actions every period, where randomization probabilities depend on actions in the precious period (memory-1) and cooperation probabilities after unilateral defection are independent of who defected.

behavior strategies that we consistently observe on average across experiments. Further, there is no evidence for heterogeneity aside from the usual 5–10% of noise players found in most experimental studies (i.e. subjects randomizing uniformly). In many ways, this represents good news for future research, as it shows that behavior is simple and robust, it confirms a central prediction, and implicitly the underlying assumptions, of recent theoretical analyses, and it poses a puzzling question for future theoretical and behavioral research: Why do subjects play the particular behavior strategies we consistently observe across this vast set of experiments? There is no obvious answer to this. In turn, our findings put a rather tight limitation on psychological interpretations of behavior. Individual behavior in various states of the repeated game is random, implying that causality between actions in one round and reactions in subsequent rounds is stochastic. In the most psychologically interesting states, succeeding unilateral defection, cooperation probabilities are around 30% and symmetric between cooperator and defector, rendering any causal statement linking action and reaction very weak.

2 Background information

Definitions The *prisoner’s dilemma* (PD) involves two players choosing whether to cooperate (c) or defect (d). In the normalized PD each player’s payoff is 1 if both cooperate and 0 if both defect. If exactly one player cooperates, the cooperating player’s payoff is $-l$ ($l > 0$) and the defecting player’s payoff is $1 + g$ ($g > 0$). This constituent game is repeated infinitely often. Assuming players are risk neutral and discount future payoffs exponentially (using factor $\delta < 1$), this game is strategically equivalent to an indefinitely repeated one that is terminated with probability $1 - \delta$ after each round. We will refer to these games jointly as repeated prisoner’s dilemma (or, *supergame*). Given $g, l > 0$, cooperation is dominated in the one-shot game but may be sustained along the path of play in subgame-perfect equilibria of the repeated PD (depending on δ).

A *strategy* σ in the repeated PD maps all finite histories to the probability of cooperation in the next round. The strategy has *memory-1* if it prescribes the same cooperation probability for any two histories that do not differ in the actions chosen in their respective last rounds. It has *memory-2* if the same holds for the respective last two rounds. We distinguish between the round-1 action and the continuation strategy, which prescribes behavior from round 2 on, and shall focus exclusively on analyzing continuation strategies. Thus, a memory-1 strategy may be denoted as $\sigma = (\sigma_{cc}, \sigma_{cd}, \sigma_{dc}, \sigma_{dd})$ corresponding to the four non-empty memory-1 histories $\{cc, cd, dc, dd\}$. We will refer to these as *states* in the following. For example, σ_{cd} , denotes the probability of cooperation when a player’s most recent action is c and her opponent’s most recent action is d . A strategy is a *pure strategy* if it prescribes degenerate cooperation probabilities after all histories ($\sigma \in \{0, 1\}$), and it is a *behavior strategy* otherwise. A *mixed strategy* randomizes on the set of pure strategies prior

to the start of the repeated game, in contrast to the behavior strategy that randomizes during the repeated game.

The data We re-analyze the exact same set of experiments reviewed in Dal Bó and Fréchette (2018). This set comprises most of the modern experiments on repeated prisoner’s dilemmas, i.e. those published since Dal Bó (2005), and consists in total of data from 12 experiments, 32 treatments, more than 1900 subjects, and almost 145,000 decisions. The set of experiments equates with the experiments listed in Table 1. A brief review is in Appendix B, but for a detailed discussion, see Dal Bó and Fréchette (2018). Thanks to its enormous size, the wide range of experiments covered (from different experimenters in various universities and various countries), and its comprehensive character with respect to the recent list of experiments on the repeated PD, this data set appears to be optimal for our purpose. In addition, by sticking exactly to the list of experiments reviewed by Dal Bó and Fréchette (2018), we can rule out the notion that data selection biases the results in favor of the hypotheses we intend to test. Finally, our analysis of continuation strategies complements the analysis by Dal Bó and Fréchette (2018), who focus largely on aggregate and first round cooperation rates as a function of treatment.

Table 1 provides an overview of behavior across experiments. It reports the average cooperation rates across experiments in each of the memory-1 states and tests for significance of difference. Besides clarifying average behavior, this allows us to test whether the surprising observation of Breitmoser (2015) that average strategies have a “semi-grim” pattern was specific to the four experiments he analyzed or can be considered a general phenomenon. A behavior strategy is called *semi-grim* if $\sigma_{cc} > \sigma_{cd} \approx \sigma_{dc} > \sigma_{dd}$, and indeed, this applies across most experiments, both for inexperienced subjects (the first half of sessions) and experienced subjects (second half of sessions). The differences between inexperienced and experienced subjects are clearly minor overall, the aggregate cooperation probabilities shift by at most five percentage points, but since it is customary to distinguish experienced and inexperienced behavior, we will consistently do so throughout this paper.

Our main hypothesis is that these semi-grim strategies reflect behavior not just on average, but also individually, as opposed to reflecting the sum of a variety of deterministic choice rules. The results of a very simple test of this hypothesis are reported in the last four columns of Table 1. These columns list the number of subjects (per experiment) that deviate significantly from randomizing 50-50 in the four memory-1 states. We focus on subjects with at least five observations per state, which suffices to trigger significance in two-sided Fisher tests if subjects play a pure strategy. The results are fairly revealing: In state (c, d) , i.e. after unilateral defection of the opponent, all standard pure strategies (except “always cooperate”) agree on the (pure) prediction that one should defect. This state is unique with respect to the unanimity of the prediction. For this state, however, we find the lowest number of subjects significantly deviating from randomizing 50-50—only around a quarter of the subjects do so, putting a rather

Table 1: Few subjects play pure strategies and assuming pure strategies yields a striking bias even in large mixture models

Experiment	Actual cooperation rates				Best-fitting rates assuming pure str.				Number of subjects not randomizing 50-50						
	$\hat{\sigma}_{cc}$	$\hat{\sigma}_{cd}$	$\hat{\sigma}_{dc}$	$\hat{\sigma}_{dd}$	$\tilde{\sigma}_{cc}$	$\tilde{\sigma}_{cd}$	$\tilde{\sigma}_{dc}$	$\tilde{\sigma}_{dd}$	(c, c)	(c, d)	(d, c)	(d, d)			
First halves per session															
<i>Aoyagi and Frechette (2009)</i>	0.917	>>	0.45	≈	0.408	≈	0.336	0.818 ⁻	0.449	0.419	0.336	32/38	1/23	3/20	7/21
<i>Blonski et al. (2011)</i>	0.89	>>	0.279	≈	0.193	>>	0.034	0.891	0.18 ⁻	0.184	0.053	13/17	1/5	3/3	124/135
<i>Bruttel and Kamecke (2012)</i>	0.91	>>	0.286	≈	0.228	>>	0.08	0.882	0.158 ⁻⁻	0.18	0.121	12/18	6/23	8/21	32/36
<i>Dal Bó (2005)</i>	0.922	>>	0.212	<	0.342	>>	0.089	0.896	0.16 ⁻	0.346	0.108	13/13	0/3	2/2	42/54
<i>Dal Bó and Fréchet (2011)</i>	0.951	>>	0.334	≈	0.331	>>	0.063	0.884 ⁻	0.196 ⁻⁻	0.318	0.108	94/106	28/117	51/128	218/253
<i>Dal Bó and Fréchet (2015)</i>	0.94	>>	0.297	≈	0.335	>>	0.057	0.891	0.172 ⁻⁻	0.317	0.09	216/243	37/137	62/147	404/474
<i>Dreber et al. (2008)</i>	0.904	>>	0.217	≈	0.213	>>	0.036	0.915	0.081 ⁻⁻	0.199	0.084	15/25	3/19	12/18	45/48
<i>Duffy and Ochs (2009)</i>	0.904	>>	0.301	≈	0.33	>>	0.111	0.863	0.239 ⁻	0.342	0.15	43/57	4/25	10/24	61/82
<i>Fréchet and Yuksel (2017)</i>	0.943	>>	0.141	≈	0.266	≈	0.091	0.918	0.099	0.274	0.085	21/28	0/0	2/2	5/8
<i>Fudenberg et al. (2012)</i>	0.982	>>	0.4	≈	0.427	>>	0.066	0.946	0.295 ⁻⁻	0.404	0.06	38/43	1/6	5/11	20/25
<i>Kagel and Schley (2013)</i>	0.935	>>	0.263	≈	0.295	>>	0.051	0.906	0.149 ⁻⁻	0.294	0.086	71/81	20/71	32/60	98/111
<i>Sherstyuk et al. (2013)</i>	0.945	>>	0.328	≈	0.371	>>	0.117	0.868 ⁻	0.201 ⁻⁻	0.377	0.125	37/44	10/36	12/34	41/52
Pooled	0.938	>>	0.304	≈	0.322	>>	0.065	0.888	0.189 ⁻⁻	0.311	0.1	605/713	111/465	202/470	1097/1299
Second halves per session															
<i>Aoyagi and Frechette (2009)</i>	0.958	>>	0.398	≈	0.517	≈	0.375	0.901 ⁻	0.367	0.544	0.325	33/37	0/12	1/12	5/9
<i>Blonski et al. (2011)</i>	0.923	>>	0.287	≈	0.231	>>	0.02	0.92	0.189 ⁻	0.221	0.04	26/32	10/25	11/16	172/178
<i>Bruttel and Kamecke (2012)</i>	0.947	>>	0.221	≈	0.297	>>	0.041	0.942	0.128 ⁻	0.315	0.067	13/15	8/17	9/12	31/35
<i>Dal Bó (2005)</i>	0.92	>>	0.242	<	0.388	>>	0.064	0.914	0.193	0.386	0.1	18/27	0/3	0/1	50/65
<i>Dal Bó and Fréchet (2011)</i>	0.979	>>	0.376	≈	0.362	>>	0.041	0.957	0.235 ⁻⁻	0.364	0.073	132/137	34/89	62/100	196/215
<i>Dal Bó and Fréchet (2015)</i>	0.976	>>	0.315	<	0.402	>>	0.035	0.947	0.187 ⁻⁻	0.408	0.061	340/365	52/162	77/146	448/497
<i>Dreber et al. (2008)</i>	0.917	>>	0.128	≪	0.39	>>	0.009	0.936	0.087	0.395	0.044	14/18	6/11	6/12	41/43
<i>Duffy and Ochs (2009)</i>	0.977	>>	0.367	≈	0.391	>>	0.082	0.923 ⁻	0.205 ⁻⁻	0.37	0.085	80/87	5/35	16/43	60/68
<i>Fréchet and Yuksel (2017)</i>	0.97	>>	0.233	≈	0.398	>>	0.069	0.93	0.072 ⁻⁻	0.391	0.074	33/37	1/6	2/10	20/25
<i>Fudenberg et al. (2012)</i>	0.971	>>	0.487	≈	0.412	>>	0.083	0.942	0.43 ⁻	0.378	0.091	41/44	2/8	4/10	14/17
<i>Kagel and Schley (2013)</i>	0.966	>>	0.262	≈	0.332	>>	0.025	0.947	0.175 ⁻	0.351	0.05	87/90	16/56	30/46	91/97
<i>Sherstyuk et al. (2013)</i>	0.973	>>	0.482	≈	0.437	>>	0.078	0.919 ⁻	0.375 ⁻⁻	0.4	0.122	44/48	7/24	17/23	23/29
Pooled	0.971	>>	0.327	<	0.376	>>	0.039	0.941	0.209 ⁻⁻	0.376	0.064	861/937	141/448	235/431	1151/1278

Note: The “actual cooperation rates” are the relative frequencies estimated directly from the data. The relation signs encode bootstrapped p -values (resampling at the subject level with 10,000 repetitions) where $<$, $>$ indicate rejection of the Null of equality at $p < .05$ and \ll , \gg indicating $p < .002$. Following Wright (1992), we accommodate for the multiplicity of comparisons within data sets by adjusting p -values using the Holm-Bonferroni method (Holm, 1979). Note that all details here exactly replicate Breitmoser (2015). As a result, if a data set is considered in isolation, the .05-level indicated by “ $>$, $<$ ” is appropriate. If all 24 treatments are considered simultaneously, the corresponding Bonferroni correction requires to further reduce the threshold to $.002 \approx .05/24$, which corresponds with “ \gg , \ll ”. The “best-fitting rates assuming pure strategies” provide the cooperation probabilities explained assuming mixtures of pure strategies as usually analyzed in the literature (the five 1-memory and five 2-memory strategies usually considered). The “number of subjects not randomizing 50-50” indicates the number of subjects with cooperation rates in the various states differing significantly from 50-50 (in subject-level two-sided binomial tests), conditioning on subjects having moved at least five times in the respective state. The required level of significance is set at $p = 0.0625$ such that five observations are sufficient to trigger statistical significance if the subject plays a pure strategy.

tight bound on the number of subjects potentially playing pure strategies.

The middle set of columns serves to illustrate the basic deficiency of deterministic choice rules. Assume that subjects use deterministic choice rules. Given that the semi-grim pattern results on average, there have to be subjects that systematically cooperate after unilateral defection of opponents. They do not have to cooperate always, but fairly often to make up for the general tendency toward defection in state (c, d) predicted by the standard strategies. For example, Result 6 of Dal Bó and Fréchette (2018) states that “always defect” (AD), Grim, and tit-for-tat (TFT) are the “three strategies [that] account for most of the data”—which seems to directly contradict the observation that $\sigma_{cd} \approx \sigma_{dc}$ unless around a third of the subjects systematically cooperate in state (c, d) . Now allow for arbitrary mixtures of the ten pure strategies that had been found to be significant in previous analyses,² and let the mixture weights be adapted independently to each treatment of each experiment to optimally capture behavior. The middle set of columns compares the cooperation probabilities predicted by these best fitting mixtures (the econometric details follow in Section 3) with the actual cooperation probabilities by experiment. The quintessence is the strong bias in state (c, d) , where predicted cooperation rates are substantially below the actual cooperation rates across experiments. Subjects cooperate much more often after unilateral defection of opponents than is compatible with the notion of pure strategies. The strategies predicting at least occasional cooperation after (c, d) , such as always-cooperate and tit-for-2-tats, simply do not fit behavior of sufficiently many subjects to capture aggregate behavior. Based on this, we will obtain the range of highly significant results suggesting that subjects play behavior strategies reported below.

Related behavioral literature We keep the literature review short and focused due to the availability of the excellent recent survey by Dal Bó and Fréchette (2018). The modern experimental research on the repeated prisoner’s dilemma started with Dal Bó (2005), who criticized earlier experiments for the implemented experimental designs such as letting subjects play against computerized opponents. The first wave of experiments following Dal Bó (2005) includes Dreber et al. (2008), Duffy and Ochs (2009), Blonski et al. (2011) and Kagel and Schley (2013), and focuses on analyzing first-round and total cooperation rates. A second wave comprising Dal Bó and Fréchette (2011, 2015), Bruttel and Kamecke (2012), Camera et al. (2012), Fudenberg et al. (2012), Sherstyuk et al. (2013), Breitmoser (2015), and Fréchette and Yuksel (2017) analyzes the continuation strategies following round-1 that we also focus on. The general theme in the reported results is that initial cooperation rates depend on the strategic environment. More specifically, they show that subgame perfection of grim is necessary but not sufficient for cooperation to emerge, and that subsequent cooperation of

²These are five memory-1 strategies and five memory-2 strategies: tit-for-tat (TFT), grim, win-stay-lose-shift (WSLS), always defect (AD), always cooperate (AC), Grim2, tit-for-2-tats (TF2T), 2-tits-for-tat (2TFT), win-stay-lose-shift-2, and “T2” (i.e. punish defection for two periods, otherwise cooperate). The standard definitions are provided in Table 5 in the online appendix.

subjects depend on their opponent's actions, primarily on those in the previous round. Many of the second-wave analyses classify the strategies of individual subjects into varying pre-selected sets of pure strategies. These analyses typically conclude that the majority of subjects plays either AD, TFT, or Grim (Dal Bó and Fréchette, 2018, Result 6), with TFT being attributed population weights around 30% if the candidate set includes only pure strategies.

Breitmoser (2015) points out the misalignment of this claim with cooperation rates state-by-state. In a data set comprising four experiments he analyzes whether semi-grim behavior strategies better capture behavior than mixtures of pure memory-1 strategies, assuming each subject consistently plays a single memory-1 strategy throughout the experiment.³ The results are suggestive, but obviously inconclusive, as the data set might be fortunately selected, as behavior might be more complex than memory-1 admits, and as subjects might switch strategies as the session progresses. In this paper, we report the results of an analysis relaxing all three concerns to arrive at arguably conclusive results. The data concern was addressed above. The case for memory-2 strategies had been made by Fudenberg et al. (2012), who show that if we assume subjects play pure strategies, then there must be subjects with memory-2, based on evidence for 2TFT and "forgiving" Grim2 strategies. Similar ideas are expressed in Aoyagi and Fréchette (2009) and Bruttel and Kamecke (2012). In contrast to these studies, we relax the assumption of pure strategies in the first place, noting that behavior strategies themselves generate behavior resembling memory 2 or 3. Recently, Fréchette and Yuksel (2017) evaluate a pre-selected set of 20 memory-1 and memory-2/3 strategies and observe no evidence for memory-2/3 (in their random termination treatment). This observation does not contradict the earlier statements, as it may result due to the large set of candidate strategies considered in the first place, which amplifies standard errors, in particular if a single treatment is considered. Our analysis resolves such concerns using the data-mining techniques described in the next section.

3 Methodology and first results

Econometric approach Recall that a subject using a pure strategy acts equivalently whenever a given state is reached and she uses the same pure strategy across all supergames. A subject using a mixed strategy uses a pure strategy within supergames but randomizes over pure strategies between supergames. A subject using a behavior strategy systematically deviates from pure strategies even within supergames. These definitions provide a basis for identification, given the set of pure strategies considered in the behavioral literature, but identification is made difficult by the standard assumption that choice is stochastic. Specifically, a single deviation from a given pure strategy, over say 20 observations, is intuitively not considered sufficient evidence against pu-

³The only other study investigating a behavior strategy seems to be Fudenberg et al. (2012), who include the strategy "generous TFT" which randomizes (only) after opponent's defection.

rity of strategies. Otherwise, the case for behavior strategies would be trivial, but how can this intuition be made formally precise—in a manner that allows us to distinguish “noisy” pure, mixed, and behavior strategies?

The distinction is achieved efficiently using the Markov-switching models known from empirical finance and empirical macroeconomics in conjunction with the robust likelihood-ratio tests of Schennach and Wilhelm (2017). Markov-switching models generalize the finite-mixture and random-switching models used in previous analyses of repeated game strategies.⁴ They allow us to capture a potentially heterogeneous group of agents (in our case, subjects potentially playing different strategies), where each agent is characterized by a “state of mind” (the strategy to be played), and agents may change their states of mind over the course of time, but both states and transitions are latent and thus not directly observable. The identifying assumption is that state transitions follow a Markov process, which substantially generalizes models assuming these transition probabilities to be degenerate (finite mixture) or uniform (random switching). Given this, estimation proceeds by maximum likelihood using an EM algorithm. Model adequacy is evaluated using ICL-BIC (Biernacki et al., 2000), and model differences are evaluated using the Schennach-Wilhelm test, which captures that all models may be arbitrarily nested and misspecified. Finally, we allow for stochastic choice in the form of trembles following Harless and Camerer (1994). All details are reported in Appendix A.

Data mining pure and mixed strategies Dal Bó and Fréchette (2018) argue that the five pure memory-1 strategies, namely AC, AD, TFT, Grim, and D-TFT, capture the behavior of most subjects across conditions. We focus on these strategies initially. Yet, as pure strategies have already shown to underestimate cooperation rates in state (c, d) in Table 1, we widen the scope of pure strategies by introducing generalized versions that allow for some randomization within supergames. Using the notation introduced above, generalized AC and AD are defined as behavior strategies $(\theta^{AC}, \theta^{AC}, \theta^{AC}, \theta^{AC})$ and $(\theta^{AD}, \theta^{AD}, \theta^{AD}, \theta^{AD})$ respectively, where θ^{AC} and θ^{AD} may be different. Generalized TFT is defined as $(1, 0, \theta^{TFT}, 0)$, generalized Grim as $(1, \theta^G, \theta^G, \theta^G)$, and generalized WSLS as $(1, 0, 0, \theta^{WSLS})$, with all $\theta^* \in [0, 1]$. Given these definitions, we mine the data as follows. We illustrate all steps by referring to Table 2, which provides the results for memory-1 strategies.

First, we evaluate which mixture of pure or generalized pure strategies best captures behavior, independently for each treatment. That is, we determine for each treatment, which combination of *pure* strategies fits best, which combination of *generalized*

⁴The approach of using mixture models in order to uncover decision rules in experimental data has been established by Stahl and Wilson (1994) and El-Gamal and Grether (1995) and subsequently used in many analyses of level- k reasoning and stochastic choice, see e.g. Houser and Winter (2004) and Houser et al. (2004), to unravel individual decision rules. A special case of finite mixture modeling is the Strategy Frequency Estimation Method (SFEM) employed by Dal Bó and Fréchette (2011), Fudenberg et al. (2012), Rand et al. (2015), Dal Bó and Fréchette (2015), Fréchette and Yuksel (2017).

pure strategies fits best, and which of the combinations fits best overall. For simplicity, we assume the best combination always contains at least TFT, AD, and Grim, as those were reported to capture players’ behavior in previous studies. This chooses the best out of the eight most promising memory-1 mixtures, for each of the 32 treatments independently. In total, we thus evaluate 8^{32} models per level of experience. Finally, we do all of this separately for three “switching models” designed to capture changes between supergames: "No Switching" (pure strategy), "Random Switching" (mixed strategy), and "Markov Switching" (where strategy choice between supergames follows a Markov process).

The results for each of the three switching models are reported in the first three columns of Table 2. For sake of readability, we aggregate ICL-BICs by experiment.⁵ This initial analysis suggests that inexperienced subjects switch strategies randomly between supergames, while experienced subjects then play constant strategies. As a side note, we obtain strong evidence that subjects play the generalized pure strategies rather than the actual pure ones (Table 11 in the appendix provides details). The differences in model fit are large, amounting to more than 1000 points on the log-likelihood scale, suggesting that randomization within supergames is indeed a behavioral facet.

Second, given the best strategy combination at treatment level, we pick the best switching model by experience level, column 4 (“Best Switching”) of Table 2, and evaluate this model against the simple semi-grim model reported in the fifth column. Note that we implicitly pick the best-fitting model from 3×8^{32} models, which is chosen after estimating 276 parameters per treatment, and evaluate it against the three-parametric behavior strategy semi-grim ($\theta_1^{SG}, \theta_2^{SG}, \theta_2^{SG}, \theta_3^{SG}$). In line with the data-mining ideal, we do not account for the degrees of freedom used in the model selection process, but solely account for the 3–10 parameters of the best-fitting model that is finally used: Still, the simple behavior strategy fits (weakly) better than the mined mixture of pure strategies. Below, we will verify if this observation holds similarly after accounting for memory-2 behavior.

Third, we make a first robustness check to clarify if perhaps not all subjects stick with memory-1 semi-grim strategies. We evaluate if some subjects simply randomize uniformly after all histories of play in the sense of level-0. The results are reported in the seventh column (“SG + Noise”). Allowing for noise players turns out to improve goodness-of-fit significantly for both experience levels—in line with many results in the literature (Nagel, 1995; Stahl and Wilson, 1995; Costa-Gomes et al., 2001), around 5-10 percent of the subjects are noise players. We conduct more such checks challenging memory-1 semi-grim below.

Fourth, we evaluate the full-blown model testing, independently treatment by treatment, which mixture and switching model fits behavior best. Thus, we choose the best-fitting model from 24 models for each treatment, amounting to the enormous

⁵Treatmentwise ICL-BICs are provided in the appendix, after Table 11. Each entry in the aggregated table represents the sum of ICL-BICs of the best out of eight models for each respective treatment.

Table 2: Best mixtures of pure or generalized strategies in relation to semi-grim (ICL-BIC of the models, less is better and relation signs point toward better models)

	Best mixture of pure or generalized strategies								Best Mixture Best Switching By Treatment				
	No Switching	Random Switching	Markov Switching	Best Switching	Semi-Grim	SG + Noise							
Specification													
# Models evaluated	8 ³²	8 ³²	8 ³²	3 × 8 ³²	1	1			24 ³² ≈ 10 ⁴⁴				
# Pars estimated (by treatment)	48	48	180	276	3	4			276				
# Parameters accounted for	3–10	3–10	12-35	3–10	3	4			3–10				
First halves per session													
<i>Aoyagi and Frechette (2009)</i>	645.31	≈	646.53	≈	649.53	≈	646.53	≈	694.72	>	646.73	≈	645.31
<i>Blonski et al. (2011)</i>	677.42	≫	619.43	≪	821.29	≈	619.43	≫	549.45	≈	558.66	<	614.47
<i>Bruttel and Kamecke (2012)</i>	585.42	≈	570.56	≈	570.87	≈	570.56	≈	567.86	≈	559.05	≈	570.56
<i>Dal Bó (2005)</i>	394.44	≫	365.35	≪	393	≈	365.35	≈	358.51	≈	361.22	≈	365.35
<i>Dal Bó and Fréchet (2011)</i>	3536.73	≈	3576.34	≈	3524.77	≈	3576.34	≈	3533.99	>	3416.54	≈	3471.98
<i>Dal Bó and Fréchet (2015)</i>	5250.62	≫	5006.33	≈	5049.55	≈	5006.33	≈	4991.74	≈	4935.93	≈	4965.57
<i>Dreber et al. (2008)</i>	459.11	≈	463.01	≈	477.02	>	463.01	>	437.17	≈	435.65	≈	461.11
<i>Duffy and Ochs (2009)</i>	1047.59	≈	1053.04	≈	1049.79	≈	1053.04	≈	1090.22	>	1017.89	≈	1047.59
<i>Fréchet and Yuksel (2017)</i>	181.98	≫	161.75	<	175.59	≈	161.75	≈	161.45	≈	164.51	≈	161.75
<i>Fudenberg et al. (2012)</i>	319.45	≈	308.6	≈	320.55	≈	308.6	≈	291.43	≈	288.89	≈	308.6
<i>Kagel and Schley (2013)</i>	1761.98	≈	1780.97	>	1694.94	≈	1780.97	≈	1782.82	≈	1726.06	≈	1694.94
<i>Sherstyuk et al. (2013)</i>	865.67	≈	907.14	>	858.65	≈	907.14	≈	912.8	≈	868.89	≈	858.65
Pooled	15951.95	≫	15675.49	≪	16214.1	>	15675.49	>	15481.59	≫	15125.92	≪	15535.6
Second halves per session													
<i>Aoyagi and Frechette (2009)</i>	363.58	≈	368.23	≈	368.89	≈	363.58	≈	389.24	>	353.94	≈	363.58
<i>Blonski et al. (2011)</i>	946.2	>	912.16	≪	1107.23	≈	946.2	>	867.87	≈	868.31	≈	908.25
<i>Bruttel and Kamecke (2012)</i>	342.17	≈	358.12	≈	347.08	≈	342.17	≈	347.4	≈	342.6	≈	342.17
<i>Dal Bó (2005)</i>	461.23	≈	445.57	<	469.98	≈	461.23	>	424.44	≈	429.06	≈	442.59
<i>Dal Bó and Fréchet (2011)</i>	2737.11	<	2865.45	>	2721.88	≈	2737.11	≈	2817.31	≈	2694.4	≈	2700.09
<i>Dal Bó and Fréchet (2015)</i>	5153.82	≈	5067.01	≈	5106.57	≈	5153.82	≈	5043.81	>	4902.82	≈	4998.09
<i>Dreber et al. (2008)</i>	287.49	≈	281.99	≈	299.36	≈	287.49	≈	264.94	≈	268.86	≈	283.38
<i>Duffy and Ochs (2009)</i>	1381.01	≈	1416.71	≈	1392.49	≈	1381.01	≈	1403.03	>	1308.8	≈	1381.01
<i>Fréchet and Yuksel (2017)</i>	309.63	≈	304.7	≈	308.78	≈	309.63	≈	313.5	≈	278.74	<	304.7
<i>Fudenberg et al. (2012)</i>	373.44	≈	395.32	≈	376.62	≈	373.44	≈	380.75	≈	358.86	≈	373.44
<i>Kagel and Schley (2013)</i>	1170.12	≈	1224.37	>	1143.67	≈	1170.12	≈	1211.37	≈	1153.04	≈	1143.67
<i>Sherstyuk et al. (2013)</i>	527.09	≈	590.16	≈	567.63	≈	527.09	≈	586.72	≈	564.5	≈	527.09
Pooled	14269.01	≈	14448.15	<	14877.81	≈	14269.01	≈	14159.8	≫	13669.82	≪	14108.4

Note: Relation signs encode p -values of Schennach-Wilhelm likelihood-ratio tests where $<$, $>$ indicate rejection of the Null of equality at $p < .05$ and \ll , \gg indicating $p < .002$, which implements the Bonferroni correction of 24 simultaneous tests per hypothesis. “No Switching” assumes that subjects chooses a strategy prior to the first supergame and plays this strategy constantly for the entire half session. “Random Switching” assumes that subjects randomly chooses a strategy prior to each supergame (by i.i.d. draws), and “Markov Switching” allows that these switches follow a Markov process.

selection of the best out of 24^{32} models across all experiments. Note that such analysis without imposing consistency requirements across treatments does not yield economically useful estimates, but it provides an upper bound on the economic content of pure and generalized pure strategies. The results are reported in the seventh column (“Best Switching By Treatment”). In total, this exhaustively mined model still does not fit better than simple semi-grim strategy. It fits slightly worse than semi-grim for inexperienced subjects and slightly better for experienced ones, but all differences are insignificant. Yet, it fits significantly worse than semi-grim after controlling for the possibility that some subjects are pure noise players.

Result 1 (Memory-1). *Assuming subjects have memory-1, the upper bound of behavior that can be captured after mining for pure or generalized pure strategies is indistinguishable from behavior captured by semi-grim, and significantly lower than by semi-grim accounting for noise players.*

4 Relaxing memory-1 and homogeneity

Data mining the memory-2 specification Next we extend the set of pure strategies to capture possible interdependence of actions with choices in $t - 2$, again to evaluate the best fitting specification against memory-1 semi-grim. We allow for two alternative approaches of extending the set of memory-1 strategies to memory-2. One approach follows Fudenberg et al. (2012), who introduced lenient and resilient variants of the pure memory-1 strategies, e.g. , strategies that punish only after the second deviation or that punish for two rounds instead of one, respectively. This approach is applicable in particular to extend pure strategies, by providing a specific list of memory-2 generalizations. The other approach is parametric and is suitable in particular to extend generalized pure and behavior strategies from memory-1 to memory-2. It allows the cooperation probabilities in one round to depend on the behavior of one or both players in $t - 2$. We allow for three different formulations here: cooperation probabilities may be a function of the opponent’s choice in $t - 2$ (*TFT-Scheme*), a function of whether both players cooperated in $t - 2$ or not (*Grim-Scheme*), or a function of the entire choice profile in $t - 2$ (*General scheme*).

First, we mine for mixtures of pure strategies, based on the list of 10 strategies⁶ of Fudenberg et al. (2012). We keep this description short, as the pure strategies fit similarly poorly as in the memory-1 case. For each treatment, we determine the most adequate combination of strategies from a list of five combinations, thus providing a selection of the best of 5^{32} models overall. The resulting model still fits highly significantly worse than the selection of generalized pure strategies with memory-1 defined above, see the two right-most columns of Table 3.

⁶These strategies are TFT, Grim, AD, Grim2, TF2T, T2, 2TFT, 2PTFT as defined in Fudenberg et al. (2012) and also in Table 5 in the Online Appendix.

Table 3: What is the memory length? Comparison of 1- and 2-memory Semi-Grim, pure and generalized strategies (ICL-BIC of the models, less is better and relation signs point toward better models)

Specification	Memory-2 Generalizations of Semi-Grim			Semi-Grim	Best Mixtures of Generalized Pure Strategies			Best Pure M1 & M2							
	M2“General”	M2“Grim”	M2“TFT”		M2“TFT”	M2“Grim”	M1								
# Models evaluated	1	1	1	1	4 ³²	4 ³²	4 ³²	5 ³²							
# Pars estimated (by treatment)	12	6	6	3	48	48	32	32							
# Parameters accounted for	12	6	6	3	9–15	9–15	6–10	3–8							
First halves per session															
<i>Aoyagi and Frechette (2009)</i>	692.5	≈	690.85	≈	686.2	≈	694.72	>	649.23	≈	646.7	≈	645.31	≪	791.38
<i>Blonski et al. (2011)</i>	714	≫	601.67	≈	601.95	≫	549.45	≪	760.49	≈	767.03	≫	713.8	≈	703.1
<i>Bruttel and Kamecke (2012)</i>	572.14	≈	566.75	≈	567.58	≈	567.86	≈	577.54	≈	581.45	≈	585.42	≈	588.55
<i>Dal Bó (2005)</i>	385.61	>	367.94	≈	366.48	≈	358.51	≪	405.84	≈	402.78	≈	407.86	≈	389.08
<i>Dal Bó and Fréchette (2011)</i>	3596.64	≈	3542.28	≈	3538.64	≈	3533.99	≈	3511.26	≈	3527.68	≈	3536.73	≪	3835.75
<i>Dal Bó and Fréchette (2015)</i>	5017.27	≈	4974.8	≈	4988.94	≈	4991.74	≪	5269.9	≈	5320.57	≈	5259.64	≪	5538.37
<i>Dreber et al. (2008)</i>	464.84	>	444.11	≈	444.71	≈	437.17	≈	477.33	≈	483.11	≈	478.09	≈	462.71
<i>Duffy and Ochs (2009)</i>	1060.26	≈	1063.66	≈	1074.9	≈	1090.22	≈	1036.79	≈	1044.14	≈	1047.59	≈	1102.64
<i>Fréchette and Yuxsel (2017)</i>	174.64	≈	167.06	≈	164.75	≈	161.45	≪	182.81	≈	188.15	≈	188.5	≈	181.98
<i>Fudenberg et al. (2012)</i>	301.76	≈	293.52	≈	294.4	≈	291.43	<	319.76	≈	322.36	≈	319.45	<	366.78
<i>Kagel and Schley (2013)</i>	1746.26	≈	1749.95	≈	1753.68	≈	1782.82	>	1634.47	≈	1665.36	<	1761.98	≈	1805.95
<i>Sherstyuk et al. (2013)</i>	917.07	≈	907.95	≈	913.52	≈	912.8	≈	872.66	≈	873.21	≈	865.67	<	941.91
Pooled	16080.69	≫	15589.39	≈	15614.59	>	15481.59	≪	16083.49	≈	16207.37	≈	16077.95	≪	16858.86
Second halves per session															
<i>Aoyagi and Frechette (2009)</i>	396.32	≈	391.42	>	387.48	≈	389.24	≈	368.09	≈	365.6	≈	363.58	≪	484.41
<i>Blonski et al. (2011)</i>	1012.48	≫	919.29	≈	922.48	≫	867.87	<	1005.42	≈	1020.36	>	992.44	≈	1055.95
<i>Bruttel and Kamecke (2012)</i>	333.51	≈	337.12	≈	329.73	≈	347.4	≈	324.09	≈	336.35	≈	344.88	≈	316.38
<i>Dal Bó (2005)</i>	449.03	≈	434.38	≈	433.82	≈	424.44	<	451.34	<	471.65	≈	475.11	≈	463.54
<i>Dal Bó and Fréchette (2011)</i>	2854.52	≈	2801.46	≈	2800.71	≈	2817.31	≈	2664.52	≈	2690.03	≈	2737.11	<	2885.43
<i>Dal Bó and Fréchette (2015)</i>	5006.3	≈	5013.49	≈	5012.99	≈	5043.81	≈	5100.98	<	5202.24	≈	5164.78	≪	5577.55
<i>Dreber et al. (2008)</i>	272.94	≈	258.88	≈	253.47	≈	264.94	≈	287.55	≈	288.95	≈	295.06	≈	287.58
<i>Duffy and Ochs (2009)</i>	1375.43	≈	1367.68	≈	1389.92	≈	1403.03	≈	1348.96	≈	1367.92	≈	1381.01	≪	1617.77
<i>Fréchette and Yuxsel (2017)</i>	308.21	≈	304.2	≈	306.93	≈	313.5	≈	311.04	≈	311.99	≈	309.63	≪	356.11
<i>Fudenberg et al. (2012)</i>	384.37	≈	382.32	≈	378.59	≈	380.75	≈	364.71	≈	359.37	≈	373.44	<	447.19
<i>Kagel and Schley (2013)</i>	1204.38	≈	1202.61	≈	1197.19	≈	1211.37	>	1052.18	≈	1066.23	≪	1170.12	≈	1169.31
<i>Sherstyuk et al. (2013)</i>	598.79	≈	590.65	≈	591.38	≈	586.72	>	510.29	≈	523.29	≈	527.09	≈	583.8
Pooled	14633.97	≫	14222.35	≈	14223.53	≈	14159.8	≈	14152.67	≪	14384.17	≈	14387.48	≪	15402.36

Note: The main body contains ICL-BICs aggregated at paper level. Relation signs and p -values are exactly as above, see Table 2. “M2” (“M1”) denotes strategies, whose actions may depend on actions in $t - 2$ and $t - 1$ ($t - 1$ only). The supplements “General”, “TFT”, “Grim” indicate whether parameters of behavior strategies may depend on: all four possible histories in $t - 2$ (M2 “General”), whether the opponent cooperated in $t - 2$ (M2 “TFT”), or whether there was joint cooperation in $t - 2$ (M2 “Grim”). Pure M2 strategies do not have such free parameters. Columns 1-3 contain one memory-2 version of semi-grim each. Column 4 is memory-1 semi-grim. Columns 5-7 are memory-2 and memory-1 versions of generalized prototypical strategies. The last column contains the best fitting combinations of a set of pure memory-1 and memory-2 strategies from the literature (TFT, Grim, AD, Grim2, TF2T, T2, 2TFT, 2PTFT) for definitions see Table 5 in the Online Appendix.

Second, we evaluate extensions of the generalized pure strategies to memory-2, using both the TFT-scheme and the Grim-scheme as defined above. That is, the strategy parameters $(\theta^{AC}, \theta^{AD}, \theta^{TFT}, \theta^G, \theta^{WSLS})$ are allowed to be functions of the opponent's action in $t - 2$ (TFT-scheme) or of whether both players cooperated in $t - 2$ (Grim-scheme). For each treatment, we determine the best of four combinations exactly as in the memory-1 case, thus determining the best of 4^{32} specifications overall. The results are reported in three columns headlined "Best Mixtures of Generalized Pure Strategies" in Table 3. In general, the TFT-scheme seems to fit better than the Grim-scheme, and overall, allowing for memory-2 helps to better capture behavior only if subjects are experienced. This suggests that behavior becomes more nuanced as subjects gain experience, but the resulting model does still not improve on the simple semi-grim strategy. Note that this obtains despite the enormous number of free parameters used in the estimation, 48 by treatment instead of 3—the pure and generalized pure strategies plainly have the wrong structure, for the reasons outlined already in section 2.

Third, we proceed similarly for evaluating extensions of semi-grim toward memory-2. That is, we compare the simple three-parametric memory-1 version with three generalizations to memory-2. As above, the TFT-scheme allows the cooperation probabilities to be functions of the opponent's action in $t - 2$, the Grim-scheme allows them to be functions of whether both subjects cooperated in $t - 2$, and the General scheme of all four possible states in $t - 2$. The results are again clear-cut: None of the memory-2 extensions improves on describing behavior by the simple memory-1 semi-grim strategy. Indeed, the finer the memory-2 ramifications, the worse the model adequacy (after accounting for the additional degrees of freedom). These results are additionally compatible with a result of Breitmöser (2015) who verified the Markov assumption by testing whether subjects systematically deviate from memory-1 strategies after particular histories in memory-2. Let us summarize this as follows.

Result 2. *Model adequacy does not improve by equipping subjects with memory-2, neither for (generalizations of) pure strategies nor for semi-grim.*

As a corollary, Result 1 on memory-1 behavior appears to be a generally adequate statement. In order to further verify its adequacy, we analyze if behavior of residual subjects is better captured by non-trivial strategies than by the simple Noise strategy from above. That is, are subject pools in the repeated prisoner's dilemma significantly heterogeneous beyond the existence of noise players? Data-mining exactly as above, we evaluate a wide range of specifications to capture behavior of the residual players, in order to obtain an upper bound of model adequacy for classes of specifications, and evaluate the best possible combination of specifications across treatments against the simple alternative that these residual players are noise players.

First, we evaluate the adequacy of pure or generalized pure strategies to complement semi-grim. Specifically, we evaluate finite-mixture models combining semi-grim with each of the five pure memory-1 strategies (AD, AC, TFT, WSLS, and Grim) or each of the four generalized pure strategies (generalized AD and generalized AC are

Table 4: Heterogeneity beyond semi-grim: Is there are significant secondary subject type?

	Mixtures of SG with pure/generalized pure				SG + Noise	Mixing Semi-Grim with							
	Best Pure	Best Gen	Pure or Gen	SG		SG-M2“Grim”	SG-M2“TFT”						
Specification													
# Models evaluated	5 ³²	4 ³²	9 ³²		1	1	1	1			1		
# Pars estimated (by treatment)	20	20	40		4	7	10	10			10		
# Parameters accounted for	4	5	4–5		4	7	10	10			10		
First halves per session													
<i>Aoyagi and Frechette (2009)</i>	681.99	≈	649.86	≈	649.86	≈	646.73	≈	650.83	≈	651.47	≈	656.3
<i>Blonski et al. (2011)</i>	576.49	≈	600.48	≫	568.43	≈	558.66	≪	642.5	≪	711.05	≈	717.67
<i>Bruttel and Kamecke (2012)</i>	562.35	≈	559.64	≈	559.64	≈	559.05	≈	553.59	≈	555.66	≈	553.75
<i>Dal Bó (2005)</i>	367.32	≈	363.9	≈	363.9	≈	361.22	<	379.63	≈	388.26	≈	390.43
<i>Dal Bó and Fréchette (2011)</i>	3413.18	≈	3379.17	≈	3368.28	≈	3416.54	≈	3335.4	≈	3386.94	≈	3375.88
<i>Dal Bó and Fréchette (2015)</i>	4950.44	≈	4920.88	≈	4905.86	≈	4935.93	≈	4958.73	≈	4990.23	≈	5005.61
<i>Dreber et al. (2008)</i>	432.72	≈	439.1	≈	432.72	≈	435.65	≈	438.35	≈	443.27	≈	441.52
<i>Duffy and Ochs (2009)</i>	1074.22	>	1017.65	≈	1017.65	≈	1017.89	≈	1016.91	≈	1023.24	≈	1016.46
<i>Fréchette and Yuksel (2017)</i>	165.36	≈	167.94	≈	165.36	≈	164.51	<	179.6	≈	181.01	≈	181.57
<i>Fudenberg et al. (2012)</i>	292.47	≈	291.68	≈	291.68	≈	288.89	≈	287.31	<	308.85	≈	308.99
<i>Kagel and Schley (2013)</i>	1685.56	≈	1690.55	≈	1685.56	≈	1726.06	>	1677.99	≈	1655.61	≈	1639.86
<i>Sherstyuk et al. (2013)</i>	887.79	≈	866.46	≈	866.46	≈	868.89	≈	868.11	≈	877.05	≈	877.26
Pooled	15272.25	≈	15166.16	≈	15078.27	≈	15125.92	≈	15244.28	≪	15537.4	≈	15530.04
Second halves per session													
<i>Aoyagi and Frechette (2009)</i>	392.87	>	357.21	≈	357.21	>	353.94	≈	358.14	≈	358.14	≈	359.37
<i>Blonski et al. (2011)</i>	861.53	≈	876.85	≫	855.06	≈	868.31	≈	906.87	≪	978.33	≈	982.05
<i>Bruttel and Kamecke (2012)</i>	336.68	>	323.04	≈	323.04	≈	342.6	≈	322.87	≈	314.07	≈	308.75
<i>Dal Bó (2005)</i>	431.34	≈	435.55	≈	431.34	≈	429.06	≪	455.9	≪	478.45	≈	467.72
<i>Dal Bó and Fréchette (2011)</i>	2597.41	≈	2543.05	≈	2537.58	≪	2694.4	>	2545.23	≈	2508.92	≈	2556.33
<i>Dal Bó and Fréchette (2015)</i>	4946.41	≈	4884.99	≈	4868.35	≈	4902.82	≈	4879.57	≈	4882.52	≈	4958.15
<i>Dreber et al. (2008)</i>	262.93	≈	258.73	≈	256.77	≈	268.86	≈	268.1	≈	252.76	≈	254.09
<i>Duffy and Ochs (2009)</i>	1379.73	>	1294.4	≈	1294.4	≈	1308.8	≈	1284.59	≈	1278.88	≈	1290.71
<i>Fréchette and Yuksel (2017)</i>	311.64	≈	283.31	≈	283.31	≈	278.74	≈	277.03	≈	278.37	≈	274.84
<i>Fudenberg et al. (2012)</i>	366.79	≈	364.37	≈	364.37	≈	358.86	≈	348.86	≈	351.64	≈	352.94
<i>Kagel and Schley (2013)</i>	1126.89	≈	1154.97	≈	1126.89	≈	1153.04	≈	1111.19	≈	1086.76	≈	1084.65
<i>Sherstyuk et al. (2013)</i>	539.35	≈	565.48	≈	539.35	≈	564.5	≈	544.85	≈	543.88	≈	538.39
Pooled	13735.95	≈	13560.79	≈	13438.59	<	13669.82	≈	13558.53	≈	13677.46	≈	13792.73

Note: The main body contains ICL-BICs aggregated at paper level. The first three columns pick the best fitting combinations of SG with pure, generalized, and generalized or pure memory-1 strategies. The last three columns combine memory-1 SG with a second memory-1 SG strategy and two different versions of a memory-2 SG strategy respectively. As above, the relation signs encode p -values of Schennach-Wilhelm likelihood-ratio tests where $<$, $>$ indicate rejection of the Null of equality at $p < .05$ and \ll , \gg indicating $p < .002$, which implements the Bonferroni correction of 24 simultaneous tests per hypothesis.

equivalent here). We then determine which of these models best captures behavior, treatment by treatment, to identify upper bounds of model adequacy for augmenting semi-grim with pure or generalized pure strategies. Table 4 presents the results. The first column “Best Pure” refers to the best treatment-wise combinations of semi-grim with pure strategies, “Best Gen” refers to combinations with generalized pure strategies, and “Pure or Gen” refers to the best combination with pure or generalized pure strategies. Doing so treatment by treatment, this yields a selection from 9^{32} models. Comparing aggregate ICL-BICs of the mined model to those containing noise players as complement to semi-grim indicates that there is little scope for improvement in this direction. There is no significant difference in model fit if subjects are inexperienced, and there is only one experiment for which the combination of semi-grim with a (generalized) pure strategy performs significantly better than SG+Noise. For this one experiment the best mixture is a combination of semi-grim with generalized AD/AC, which is very similar to noise: It represents unconditional randomization between cooperation and defection but allows average cooperation rates to deviate from 0.5 (at the expense of a degree of freedom).

Second, we evaluate if a residual component is better captured using a memory-1 or memory-2 semi-grim (SG) strategy instead of noise. As memory-2 implementations we evaluate the TFT- and Grim-schemes introduced above. The overall picture can be summarized succinctly: Overall, none of these models improves on the null that the residual players simply randomize fifty-fifty. The results are presented in the three right-most columns of Table 4. The aggregate ICL-BICs of the two models combining semi-grim with memory-2 semi-grim do not significantly differ from each other, both perform slightly worse than the simpler model combining memory-1 SG with another memory-1 SG (“SG+SG”), which in turn does not perform better than the SG+Noise model. Since all of these semi-grim models contain uniform randomization (noise) as special case, we conclude that the additional structure provided by these models does not improve the goodness-of-fit sufficiently to warrant the free parameters used in the process. Noise is an adequate description for the residual players.

Result 3. *There is significant evidence for subject heterogeneity, but no single model of behavior significantly improves on the hypothesis that the residual component besides semi-grim subjects consists simply of “noise players” randomizing fifty-fifty.*

5 Conclusion

Do we play pure, mixed, or behavior strategies in the repeated prisoner’s dilemma? What is the memory length of these strategies and are subjects heterogeneous? After data mining for the best-fitting mixtures of pure strategies and switching models, which puts pure and mixed strategies at an enormous advantage, we still have to reject the hypothesis that subjects play mixtures of deterministic strategies in favor of the much simpler explanation that subjects actually play the strategies that we observe

consistently on average. Even in our large data set, we find no evidence for the use of memory-2 strategies nor for subject heterogeneity beyond the existence of noise players. That is, aggregating across 12 experiments and 145,000 decisions, subjects play memory-1 behavior strategies highly consistently. These strategies depend to some extent on the treatment parameters, see Table 1, and understanding this interdependence is thus established as the key to understanding cooperation in long-run relationships.

This suggests that subjects decide “just in time” whether to retaliate or be lenient. From the position of an observer, this is best described as a random Markov process (a memory-1 behavior strategy), as opposed to a complex set of deterministic rules. From the position of the decision maker herself, any decision may appear to be perfectly deterministic, but this discrepancy has been labeled the “illusion of conscious will” (Wegner, 2004) as any decision is preceded by unconscious cerebral activity feeding us thoughts (Libet, 1985). From a game-theoretic perspective, individual decisions are not sufficiently consistent to organize them by means of deterministic rules, but an illusion of conscious will would be reminiscent of purification in the sense of Harsanyi (1973).

Our results clarify which statistics to analyze—memory-1 cooperation rates—in experimental studies of long-run interactions and underscore the practical relevance of the current theoretical literature on imperfect private monitoring, which actually predicts that players would play simple behavior strategies. The strategies observed in our analysis are very specific instances of belief-free equilibria, however, played consistently across a wide range of experiments and treatments, which is not an obvious implication of the theory as it is currently developed. The hypothesis that equilibria are additionally robust to utility perturbations (Breitmoser, 2015) may be worth investigating, but it should be noted that the cooperation probabilities exhibit semi-grim patterns even if belief-free equilibria sustaining cooperation do not exist. Analogical reasoning in the sense of Samuelson (2001) may help explain these “anomalies”, but more generally, the consistency of the semi-grim pattern in cooperation probabilities seems to represent a non-trivial challenge for future research—despite the already impressive body of theoretical work on repeated games.

References

- Aoyagi, M. and Frechette, G. (2009). Collusion as public monitoring becomes noisy: Experimental evidence. *Journal of Economic theory*, 144(3):1135–1165.
- Bajari, P. and Hortacsu, A. (2005). Are structural estimates of auction models reasonable? evidence from experimental data. *Journal of Political Economy*, 113(4):703–741.
- Biernacki, C., Celeux, G., and Govaert, G. (2000). Assessing a mixture model for clustering with the integrated completed likelihood. *IEEE transactions on pattern analysis and machine intelligence*, 22(7):719–725.
- Bilmes, J. A. et al. (1998). A gentle tutorial of the em algorithm and its application to parameter estimation for gaussian mixture and hidden markov models. *International Computer Science Institute*, 4(510):126.
- Blonski, M., Ockenfels, P., and Spagnolo, G. (2011). Equilibrium selection in the repeated prisoner’s dilemma: Axiomatic approach and experimental evidence. *American Economic Journal: Microeconomics*, 3(3):164–192.
- Breitmoser, Y. (2015). Cooperation, but no reciprocity: Individual strategies in the repeated prisoner’s dilemma. *American Economic Review*, 105(9):2882–2910.
- Brunner, C., Camerer, C. F., and Goeree, J. K. (2011). Stationary concepts for experimental 2 x 2 games: Comment. *American Economic Review*, 101(2):1029–40.
- Bruttel, L. and Kamecke, U. (2012). Infinity in the lab. how do people play repeated games? *Theory and Decision*, 72(2):205–219.
- Camera, G., Casari, M., and Bigoni, M. (2012). Cooperative strategies in anonymous economies: An experiment. *Games and Economic Behavior*, 75(2):570–586.
- Charness, G. and Rabin, M. (2002). Understanding social preferences with simple tests. *The Quarterly Journal of Economics*, 117(3):817–869.
- Costa-Gomes, M., Crawford, V. P., and Broseta, B. (2001). Cognition and behavior in normal-form games: An experimental study. *Econometrica*, 69(5):1193–1235.
- Dal Bó, P. (2005). Cooperation under the shadow of the future: experimental evidence from infinitely repeated games. *American Economic Review*, 95(5):1591–1604.
- Dal Bó, P. and Fréchet, G. R. (2011). The evolution of cooperation in infinitely repeated games: Experimental evidence. *American Economic Review*, 101(1):411–429.

- Dal Bó, P. and Fréchette, G. R. (2015). Strategy choice in the infinitely repeated prisoners' dilemma. *Working paper*, Available at SSRN: <https://ssrn.com/abstract=2292390> or <http://dx.doi.org/10.2139/ssrn.2292390>.
- Dal Bó, P. and Fréchette, G. R. (2018). On the determinants of cooperation in infinitely repeated games: A survey. *Journal of Economic Literature*, 56(1):60–114.
- Dreber, A., Rand, D. G., Fudenberg, D., and Nowak, M. A. (2008). Winners don't punish. *Nature*, 452(7185):348–351.
- Duffy, J. and Ochs, J. (2009). Cooperative behavior and the frequency of social interaction. *Games and Economic Behavior*, 66(2):785–812.
- El-Gamal, M. A. and Grether, D. M. (1995). Are people bayesian? uncovering behavioral strategies. *Journal of the American statistical Association*, 90(432):1137–1145.
- Ely, J. C., Hörner, J., and Olszewski, W. (2005). Belief-free equilibria in repeated games. *Econometrica*, 73(2):377–415.
- Ely, J. C. and Välimäki, J. (2002). A robust folk theorem for the prisoner's dilemma. *Journal of Economic Theory*, 102(1):84–105.
- Eyster, E. and Rabin, M. (2005). Cursed equilibrium. *Econometrica*, 73(5):1623–1672.
- Fehr, E. and Schmidt, K. M. (1999). A theory of fairness, competition, and cooperation. *The Quarterly Journal of Economics*, 114(3):817–868.
- Fréchette, G. R. and Yuksel, S. (2017). Infinitely repeated games in the laboratory: Four perspectives on discounting and random termination. *Experimental Economics*, 20(2):279–308.
- Frühwirth-Schnatter, S. (2006). *Finite mixture and Markov switching models*. Springer Science & Business Media.
- Fudenberg, D. and Levine, D. K. (1993). Self-confirming equilibrium. *Econometrica*, 61(3):523–545.
- Fudenberg, D. and Maskin, E. S. (1986). Folk theorem for repeated games with discounting or with incomplete information. *Econometrica*, 54(3):533–554.
- Fudenberg, D., Rand, D. G., and Dreber, A. (2012). Slow to anger and fast to forgive: Cooperation in an uncertain world. *American Economic Review*, 102(2):720–749.
- Harless, D. W. and Camerer, C. F. (1994). The predictive utility of generalized expected utility theories. *Econometrica*, 62(6):1251–1289.

- Harsanyi, J. C. (1973). Games with randomly disturbed payoffs: A new rationale for mixed-strategy equilibrium points. *International Journal of Game Theory*, 2(1):1–23.
- Holm, S. (1979). A simple sequentially rejective multiple test procedure. *Scandinavian Journal of Statistics*, 6(2):65–70.
- Houser, D., Keane, M., and McCabe, K. (2004). Behavior in a dynamic decision problem: An analysis of experimental evidence using a bayesian type classification algorithm. *Econometrica*, 72(3):781–822.
- Houser, D. and Winter, J. (2004). How do behavioral assumptions affect structural inference? evidence from a laboratory experiment. *Journal of Business & Economic Statistics*, 22(1):64–79.
- Kagel, J. H. and Schley, D. R. (2013). How economic rewards affect cooperation reconsidered. *Economics Letters*, 121(1):124–127.
- Kandori, M. (2002). Introduction to repeated games with private monitoring. *Journal of Economic Theory*, 102(1):1–15.
- Kandori, M. (2011). Weakly belief-free equilibria in repeated games with private monitoring. *Econometrica*, 79(3):877–892.
- Libet, B. (1985). Unconscious cerebral initiative and the role of conscious will in voluntary action. *Behavioral and brain sciences*, 8(4):529–539.
- Mailath, G. J. and Samuelson, L. (2006). *Repeated Games and Reputations: Long-Run Relationships*. Oxford University Press.
- Nagel, R. (1995). Unraveling in guessing games: An experimental study. *American Economic Review*, 85(5):1313–1326.
- Rand, D. G., Fudenberg, D., and Dreber, A. (2015). It’s the thought that counts: The role of intentions in noisy repeated games. *Journal of Economic Behavior & Organization*, 116:481–499.
- Samuelson, L. (2001). Analogies, adaptation, and anomalies. *Journal of Economic Theory*, 97(2):320–366.
- Schennach, S. M. and Wilhelm, D. (2017). A simple parametric model selection test. *Journal of the American Statistical Association*, pages 1–12.
- Sherstyuk, K., Tarui, N., and Saijo, T. (2013). Payment schemes in infinite-horizon experimental games. *Experimental Economics*, 16(1):125–153.
- Stahl, D. O. and Wilson, P. W. (1994). Experimental evidence on players’ models of other players. *Journal of Economic Behavior & Organization*, 25(3):309–327.

- Stahl, D. O. and Wilson, P. W. (1995). On players' models of other players: Theory and experimental evidence. *Games and Economic Behavior*, 10(1):218–254.
- Wegner, D. M. (2004). Précis of the illusion of conscious will. *Behavioral and Brain Sciences*, 27(5):649–659.
- Weizsäcker, G. (2003). Ignoring the rationality of others: evidence from experimental normal-form games. *Games and Economic Behavior*, 44(1):145–171.
- Wright, S. P. (1992). Adjusted p-values for simultaneous inference. *Biometrics*, 48:1005–1013.

Online appendix

God does not play dice, but do we?

Appendix A: Markov-Switching Models

The Markov-switching model builds on the simpler and more restrictive finite mixture model, which has been established in the experimental literature by Stahl and Wilson (1994) and can be used to empirically identify a finite number K of strategies with parameter vectors θ_k . The log-likelihood function to be maximized for the finite mixture model is

$$\ln \mathcal{L}(\theta, \rho | O) = \log \left(\prod_{s \in S} p(o_s | \theta, \rho) \right) = \sum_{s \in S} \ln \sum_{k \in K} \rho_k p_k(o_s | \theta_k), \quad (1)$$

with observations O , ρ_k denoting the relative frequency of strategy k , and $p_k(o_s | \theta_k)$ denoting the probability that player s chooses action o_s given he plays strategy k ⁷.

A way to model regime switches is to replace the implicit latent indicator variable in finite mixture models (indicating the discrete types) with a hidden Markov chain (Frühwirth-Schnatter, 2006). The central assumption characterizing the learning process in Markov models is that the type of a player (or its strategy in our context) in the next period can only depend on its type in this period. More precisely if k_t is the type in period t then: $Pr(k_{t+1} | k_t, k_{t-1}, k_{t-2}, \dots, k_1) = Pr(k_{t+1} | k_t)$, where the type is hidden and cannot be observed directly.⁸ What we do observe is the action o_t , which in turn depends on the type k_t in t only: $Pr(o_t | k_t, o_{t-1}, k_{t-1}, \dots, k_1, o_1) = Pr(o_t | k_t)$ (c.f. Bilmes et al. (1998)). It implies that transitions between states are independent of time t . This assumption might be quite restrictive. For example if we want to assume that the probability of switching to a new strategy is more likely later in the game than at the beginning. Nevertheless, we can use memory-2 or memory-3 strategies if we define the state ω as a history of more than one past outcome and condition the strategy on this history of outcomes. Moreover, it is possible to interact time dependent components with switching probabilities.

Let K_t denote the state at time $t \in 1, 2, \dots, T$ and $\sigma_{kk'} = Pr(K_{t+1} = k' | K_t = k)$ define the transition probability from k to k' which is independent from t , as pointed. So σ is a $(K \times K)$ transition matrix containing transition probabilities for every pair of states, where all entries are positive and each row sums up to 1. Moreover, the state-paths are denoted by $\kappa \in K^T$ with $Pr(\kappa)$ conditional on initial weights ρ and transition

⁷For the memory-1 case $p_k(o_s | \theta_k) = \prod_t (\sigma_{\omega_{s,t}}(k))^{o_{s,t}} (1 - \sigma_{\omega_{s,t}}(k))^{1-o_{s,t}}$ with strategy $\sigma_{\omega_{s,t}}(k)^{o_{s,t}}$ for state $\omega_{s,t}(k)^{1-o_{s,t}}$.

⁸Therefore also known as the Hidden Markov Model (HMM).

probabilities σ . The probability of observing $o_{s,t}$ conditional on subject s being type k in this period is $Pr(o_{s,t}|\theta, k)$. The likelihood function is:

$$\ln \mathcal{L}(\rho, \sigma, \theta|O) = \sum_{s \in S} \ln \sum_{\kappa \in K^T} Pr(\kappa) \prod_{t \leq T} Pr(o_{s,t}|\theta, \kappa(t), t) \quad (2)$$

Due to the introduction of the transition matrix σ the number of parameters to be estimated increases dramatically. Moreover, with a naive estimation approach we would have to consider all possible state paths and be very time consuming. Therefore we choose to apply a backward-forward algorithm to calculate posteriors for estimation with the expectation maximization (EM) algorithm.

The idea of the EM-algorithm is to conduct two steps the E-step and the M-step iteratively. This way we split up every optimization step into many steps which simplifies complexity and consequently speeds up computations. In the E-step we evaluate the conditional expectation of the log-likelihood given our data O and the current parameter vector and then maximize over a reduced set of free parameters in the M-Step. The number of possible types k is pre-defined as well as the structure of their mixing parameters θ_k .

In the E-step we need to compute for all subjects for all time periods the posterior probability of component inclusion (being a specific type) and the probability to switch between two types. An efficient way to calculate those posterior probabilities is to built up on the backward-forward. First, we have the forward procedure, where we define the (joint) probability of observing the partial sequence o_{s1}, \dots, o_{st} and ending up with type k at time t :

$$\alpha_{sk}(t) = Pr(O_{s1} = o_{s1}, \dots, O_{st} = o_{st}, K_t = k) \quad (3)$$

Recursively, we can then define:

1. $\alpha_{sk}(1) = \rho_k Pr(o_{s1}|\theta, k)$ (4)
2. $\alpha_{sk'}(t+1) = \left[\sum_k \alpha_{sk}(t) \sigma_{kk'} \right] Pr(o_{st+1}|\theta, k')$
3. $Pr(o_s) = \sum_{k \in K} \alpha_{sk}(T)$

Second, for the backward procedure we define the probability of ending in the partial sequence o_{st+1}, \dots, o_{sT} given that we have started at type k at time t .

$$\beta_{sk}(t) = Pr(O_{st+1} = o_{st+1}, \dots, O_T = o_T | K_t = k) \quad (5)$$

Again we can define $\beta_{sk}(t)$ efficiently (Bilmes et al., 1998)

1. $\beta_{sk}(T) = 1$ (6)
2. $\beta_{sk}(t) = \sum_{k' \in K} \sigma_{kk'} Pr(o_{st+1} | \theta, k') \beta_{sk'}(t+1)$
3. $Pr(o_s) = \sum_{k \in K} \beta_{sk}(1) \rho_k Pr(o_{s1} | \theta, k)$

We then take advantage of the fact that the unconditional probability $Pr(o_s)$ can be defined using $\alpha_{sk}(t)$ or $\beta_{sk}(t)$ to calculate the posterior probabilities γ_{sk} and $\zeta_{skk'}$. The former is the conditional probability of being type k at time t given observations o_s :

$$\gamma_{sk}(t) = Pr(K_t = k | o_s) = \frac{Pr(o_s, K_t = k)}{Pr(o_s)} = \frac{Pr(o_s, K_t = k)}{\sum_{k' \in K} Pr(o_s, K_t = k')} = \frac{\alpha_{sk}(t) \beta_{sk}(t)}{\sum_{k' \in K} \alpha_{sk'}(t) \beta_{sk'}(t)}, \quad (7)$$

Using γ_{sk} we can define the probability of having type k in t and type k' in $t+1$ conditional on our observations as

$$\begin{aligned} \zeta_{skk'}(t) &= Pr(K_t = k, K_{t+1} = k' | o_s) = \frac{Pr(K_t = k | o_s) Pr(o_{t+1}, \dots, T, K_{t+1} = k' | K_t = k)}{Pr(o_{t+1}, \dots, T | K_t = k)} \\ &= \frac{\gamma_{sk}(t) \sigma_{kk'} Pr(o_{s,t+1} | \theta, k') \beta_{sk'}(t+1)}{\beta_{sk}(t)} \end{aligned} \quad (8)$$

(cf. Bilmes et al. (1998)).

In the M-step we maximize for each k and $t \leq T$ the function

$$LL_{kt}(\theta'_k) = \sum_{s \in S} \gamma_{sk}(t) \ln Pr(o_{st} | \theta') \rightarrow \max_{\theta'_k} ! \quad (9)$$

to yield the updated θ^{+1} when assuming that θ_{kt} does not affect the likelihood of other components k . If it does, we need to maximize $\sum_{k' \in K} LL_{kt}(\theta') \rightarrow \max_{\theta'} !$ and yield θ^{+1} .⁹

Moreover, we update ρ and σ using the posteriors from above and yield

$$\rho_k^{+1} = \frac{1}{n} \sum_{s \in S} \gamma_{sk}(1) \quad \text{and} \quad \sigma_{kk'}^{+1} = \frac{\sum_{s \in S} \sum_{t < T} \zeta_{skk'}(t)}{\sum_{s \in S} \sum_{t < T} \gamma_{sk}(t)} \quad (10)$$

The two steps are iterated until the distance between (θ, ρ, σ) and $(\theta^{+1}, \rho^{+1}, \sigma^{+1})$ gets small.

Estimation proceeds by a maximum likelihood, as usual, but as is well-known, the

⁹ θ may depend on t but does not have to.

larger the number of parameters, the larger a model's capacity to fit data—and implicitly, the larger its fallacy to overfit the data. This is conventionally captured by evaluating model adequacy based on information criteria such as BIC, which penalize for the degrees of freedom in a theoretically adequate manner. Mixture and switching models additionally contain freedom in defining the components of the subject pool, i.e. the number of subject types, which provides an additional source for overfitting aside from the number of parameters used. Following (Biernacki et al., 2000), these concerns are addressed using the information-classification likelihood Bayes-information criterion (ICL-BIC), a criterion that penalizes both model complexity and the failure of the mixture model to provide a classification in well-separated strategy clusters. We address the observation that modeling mixtures of pure, mixed, and behavior strategies induces sophisticated nesting structures, and the concern that indeed all models may be misspecified by evaluating model differences using the novel Schennach-Wilhelm likelihood ratio tests (Schennach and Wilhelm, 2017). Finally, we capture the intuition that choice is stochastic by allowing for trembles in the sense of Selten (1975): Each agent of a player picks any given action with probability no less than $\epsilon > 0$. This approach follows (Breitmoser, 2015) and, in relation to the logistic-error approach proposed by (Dal Bó and Fréchette, 2011), it has the advantage that it does not perturb choice probabilities of subjects that originally randomize already.

Appendix B: Information on the experiments re-analyzed

This section provides some background information on the experiments re-analyzed in this paper. Table 5 summarizes and defines the strategies considered by previous studies. Table 6 reviews focus and main results (in terms of identified strategies) of these studies. Table 7 reviews the numbers of subjects and observations, average parameters, and average cooperation rates for all experiments, and Table 8 provides the detailed overview by treatments.

Table 5: Pure strategies considered in behavioral analyses

Strategy	Abbreviation	Description	$(\sigma_{cc}, \sigma_{cd}, \sigma_{dc}, \sigma_{dd})^\dagger$	References
Pure Strategies Non-responsive or Memory-1				
Always Defect	AD	Always defects independent of previous outcome	(0,0,0,0)	DF11, DF15, FRD12, FY17, STS13
Always Cooperate	AC	Always cooperates independent of previous outcome	(1,1,1,1)	DF11, DF15, FRD12, FY17, STS13
Grim	G	Only cooperates after cc was last outcome	(1,0,0,0)	AF09, DF11, DF15, FRD12, FY17, STS13
Tit-for-Tat	TFT	Only plays C if opponent did last period	(1,0,1,0)	AF09, DF11, DF15, FRD12, FY17, STS13
Win-stay-Lose-Shift (aka Perfect TFT)	WSLS	Plays same strategy if it was successful, otherwise shifts	(1,0,0,1)	DF11, DF15, FRD12, FY17
False cooperator	C-to-AD	Play c in first round then AD	–	FRD12, FY17
Explorative TFT Alternator	D-TFT DC-Alt	Play d in first round then TFT Play d in first round then alternate c and d	– –	DF15, FRD12, FY17 FRD12, FY17
Trigger-with-Reversion	GwR	Like Grim but revert to cooperation after cc [‡]	(1,0,0,0)	STS13
Pure Strategies Memory-2/3				
Trigger 2 periods	T2	Player punishes defection for max. 2 periods, otherwise cooperates	$(1,0,\theta_1^*, 0)$	DF11, FY17
Tit-for-2(3)-Tats	TF2T	Defects after 2 defections	$(1,\theta_2, 1,\theta_2)$	FRD12, FY17
2-Tits-for-2-Tats	2TF2T	Defects twice after 2 defections	$(1,\theta_3,\theta_3,\theta_3)$	FRD12, FY17
2-Tits-for-1-Tats	2TFT	Defects twice after each defections	$(1,0,\theta_4,0)$	FRD12, FY17
Grim2(3/4)	G2(3)	After 2(3) defections will play D forever	$(1,\theta_5,0,0)$	FRD12, FY17, STS13
Win-stay-Lose-Shift-2	WSLS2	cooperate after (dd,dd),(cc,cc), (dd,cc) otherwise defect	–	FRD12
Explorative TF2(3)T	D-TF2(3)T	Play D in first round then TF2(3)T	–	FRD12, FY17
Explorative Grim2(3)	D-Grim2(3)	Play D in first round then Grim2(3)	–	FRD12, FY17
Behavior Strategies				
Semi-Grim**	SG	Similar to Grim but may cooperate after CD or DC.	$(1,\theta_{SG},\theta_{SG},0)$	B15
Generous TFT	GTFT	Like TFT but cooperate with prob α after CD or DD	$(1,\theta_{GT},1,\theta_{GT})$	FRD12, B15

[†] σ assigns cooperation probabilities after joint cooperation (cc), unilateral defection by opponent (cd), unilateral defection (dc), and joint defection (dd).

[‡] possible if players make mistakes.

* Vector assigning cooperation probabilities $\in \{0, 1\}$ depending on the state 2 periods ahead.

** θ_{SG} and θ_{GT} are mixing parameters $\in (0, 1)$.

References: **AF09** (Aoyagi and Fréchette, 2009), **B15** (Breitmoser, 2015), **DF11** (Dal Bó and Fréchette, 2011), **DF15** (Dal Bó and Fréchette, 2015), **FRD12** (Fudenberg et al., 2012), **FY17** (Fréchette and Yüksel, 2017), **STS13** (Sherstyuk et al., 2013)

Table 6: Overview literature

Reference	Focus	Investigation of Strategies	Strategies found
Aoyagi and Frechette (2009)	Imperfect public monitoring in PD	Mainly avg. coop. rates Mem-1, Mem-2, Threshold	Threshold strat S_0 (same threshold in state 1 & 0)
Blonski et al. (2011)	New δ^* with sucker's payoff	Avg. coop rates	–
Bruttel and Kamecke (2012)	Endgame effects	Elicitation of pure strategies discuss avg. coop. rates	– *
Camera et al. (2012)	Player's strat using finite automata	All possible pure mem-1	large share play unconditional
Dal Bó (2005)	Finitely vs infinitely repeated PD	Avg. cooperation rates	–
Dal Bó and Fréchette (2011)	Players' strategies learning model	selected mem-1 strategies SFEM	AC, AD, TFT
Dal Bó and Fréchette (2015) upd (2017)	Players' strategies	SFEM, elicitation, pure Mem-1, Mem-2 mainly preselected	AD, TFT, Grim
Dreber et al. (2008)	PD extended with punishment option	Agg. cooperation behavior	(AD, Grim, TFT)**
Duffy and Ochs (2009)	Fixed matching of players in PD	Round 1 and avg. coop. rates	–
Fréchette and Yuksel (2017)	De-coupling of expected length of game and discount factor	Avg. coop. rates, SFEM Mem-1, Mem2/3 preselected	Grim, TFT
Fudenberg et al. (2012)	Effect of noise/uncertainty on leniency	Avg. coop. rate, SFEM, 20 pure Mem-1, Mem-2(3)	AC, AD, Grim, (D)-TFT, 2TFT, Grim2
Kagel and Schley (2013)	Linear payoff transformations	Fist round coop. rates	–
Sherstyuk et al. (2013)	Payment schemes	Avg. cooperation rates, share of correctly predicted actions by selected pure strats	AD, TFT, GwR
Dal Bó and Fréchette (2018)	Determinants of cooperation (meta)	Mainly first round coop	–

* Table 4 column "Strategy" in their study indicating SG in coefficients for cd_{t-1} & cd_{t-2} .

** Reported by Fudenberg et al. (2012).

Table 7: Overview of the data sets used in the analysis

Experiment	Logistics		Parameters			Average cooperation rates							
	#Subj	#Dec	δ	g	l	$\hat{\sigma}_0$	$\hat{\sigma}_{cc}$	$\hat{\sigma}_{cd}$	$\hat{\sigma}_{dc}$	$\hat{\sigma}_{dd}$			
First halves per session													
<i>Aoyagi and Frechette (2009)</i>	38	1650	0.9	0.333	0.111	0.465	0.917	≫	0.45	≈	0.408	≈	0.336
<i>Blonski et al. (2011)</i>	200	3040	0.756	1.345	2.602	0.295	0.89	≫	0.279	≈	0.193	≫	0.034
<i>Bruttel and Kamecke (2012)</i>	36	1920	0.8	1.167	0.833	0.481	0.91	≫	0.286	≈	0.228	≫	0.08
<i>Dal Bó (2005)</i>	102	1320	0.75	0.939	1.061	0.342	0.922	≫	0.212	<	0.342	≫	0.089
<i>Dal Bó and Fréchette (2011)</i>	266	17772	0.622	1.062	1.072	0.31	0.951	≫	0.334	≈	0.331	≫	0.063
<i>Dal Bó and Fréchette (2015)</i>	672	22112	0.743	1.579	1.341	0.451	0.94	≫	0.297	≈	0.335	≫	0.057
<i>Dreber et al. (2008)</i>	50	2064	0.75	1.488	1.488	0.488	0.904	≫	0.217	≈	0.213	≫	0.036
<i>Duffy and Ochs (2009)</i>	102	3128	0.9	1	1	0.53	0.904	≫	0.301	≈	0.33	≫	0.111
<i>Fréchette and Yuksel (2017)</i>	50	800	0.75	0.4	0.4	0.737	0.943	≫	0.141	≈	0.266	≈	0.091
<i>Fudenberg et al. (2012)</i>	48	1452	0.875	0.333	0.333	0.756	0.982	≫	0.4	≈	0.427	≫	0.066
<i>Kagel and Schley (2013)</i>	114	7600	0.75	1	0.5	0.573	0.935	≫	0.263	≈	0.295	≫	0.051
<i>Sherstyuk et al. (2013)</i>	56	3052	0.75	1	0.25	0.56	0.945	≫	0.328	≈	0.371	≫	0.117
Pooled	1734	65910	0.728	1.207	1.083	0.389	0.938	≫	0.304	≈	0.322	≫	0.065
Second halves per session													
<i>Aoyagi and Frechette (2009)</i>	38	1400	0.9	0.333	0.111	0.424	0.958	≫	0.398	≈	0.517	≈	0.375
<i>Blonski et al. (2011)</i>	200	5460	0.766	1.282	2.554	0.279	0.923	≫	0.287	≈	0.231	≫	0.02
<i>Bruttel and Kamecke (2012)</i>	36	1632	0.8	1.167	0.833	0.447	0.947	≫	0.221	≈	0.297	≫	0.041
<i>Dal Bó (2005)</i>	102	1650	0.75	0.961	1.039	0.297	0.92	≫	0.242	<	0.388	≫	0.064
<i>Dal Bó and Fréchette (2011)</i>	266	19270	0.62	1.122	1.103	0.355	0.979	≫	0.376	≈	0.362	≫	0.041
<i>Dal Bó and Fréchette (2015)</i>	672	29480	0.766	1.666	1.386	0.469	0.976	≫	0.315	<	0.402	≫	0.035
<i>Dreber et al. (2008)</i>	50	1838	0.75	1.533	1.533	0.461	0.917	≫	0.128	≪	0.39	≫	0.009
<i>Duffy and Ochs (2009)</i>	102	6018	0.9	1	1	0.684	0.977	≫	0.367	≈	0.391	≫	0.082
<i>Fréchette and Yuksel (2017)</i>	50	1568	0.75	0.4	0.4	0.763	0.97	≫	0.233	≈	0.398	≫	0.069
<i>Fudenberg et al. (2012)</i>	48	1800	0.875	0.333	0.333	0.829	0.971	≫	0.487	≈	0.412	≫	0.083
<i>Kagel and Schley (2013)</i>	114	7172	0.75	1	0.5	0.704	0.966	≫	0.262	≈	0.332	≫	0.025
<i>Sherstyuk et al. (2013)</i>	56	2604	0.75	1	0.25	0.646	0.973	≫	0.482	≈	0.437	≫	0.078
Pooled	1734	79892	0.744	1.271	1.172	0.404	0.971	≫	0.327	<	0.376	≫	0.039

A-8

Note: The “average cooperation rates” are the relative frequencies estimated directly from the data. The relation signs encode bootstrapped p -values (resampling at the subject level with 10,000 repetitions) where $<$, $>$ indicate rejection of the Null of equality at $p < .05$ and \ll , \gg indicating $p < .002$. Following Wright (1992), we accommodate for the multiplicity of comparisons within data sets by adjusting p -values using the Holm-Bonferroni method (Holm, 1979). Note that all details here exactly replicate Breitmoser (2015). As a result, if a data set is considered in isolation, the .05-level indicated by “ $>$, $<$ ” is appropriate. If all 24 treatments are considered simultaneously, the corresponding Bonferroni correction requires to further reduce the threshold to $.002 \approx .05/24$, which corresponds with “ \gg , \ll ”.

Appendix C: Robustness checks for Section 3

The tables in this section replicate the tables presented in the Section 3, provide a number of robustness checks and additionally present the results treatment-by-treatment.

- Table 9 compares the “best mixtures” analyzed in the main text to the models allowing for all 1-memory types that correspond with those analyzed in the literature, e.g. Dal Bó and Fréchette (2011). Recall that the 2-memory strategies analyzed in other strings of literature are examined in Section 4. This table clarifies that focusing on the “best mixtures” for each treatment improves the goodness-of-fit of these models substantially (i.e. by at least 100 likelihood points).
- Table 11 compares the best mixtures of pure and generalized pure strategies as discussed in the main text.
- Table 13 is similar to Table 2 in the main text but focusing on the prototypical strategies in their pure form.
- Table 15 is similar to Table 2 in the main text but focusing on the prototypical strategies in their generalized form.
- Table 17 is equivalent to Table 2 in the main text.
- Table 19 examines a number of mixture models involving Semi-Grim strategies to clarify that the mixtures analyzed in Table 2 do neither misrepresent the general picture (these alternative models are not better in terms goodness-of-fit and hence need not be discussed in detail in the main text) nor present a non-robust part of the picture (these alternative models do not fit substantially worse than simple Semi-Grim either, implying that the general results on the relation of Semi-Grim to the prototypical strategies can be considered robust).

Table 9: Pure, mixed, or switching strategies? (ICL-BIC of the models, less is better and relation signs point toward better models)

	Best w/o SG			All but SG		
	No Switching	Random Switching	Markov Switching	No Switching	Random Switching	Markov Switching
Specification						
# Models evaluated	4	4	4	1	1	1
# Pars estimated (by treatment)	16	16	82	5	5	30
# Parameters accounted for	3–5	3–5	12–30	5	5	30
First halves per session						
<i>Aoyagi and Fréchette (2009)</i>	790.2	>	767.22	≈	797.89	790.2 > 767.22 ≈ 797.89
<i>Blonski et al. (2011)</i>	674.3	≫	626.71	≪	843	719.21 ≻ 674.3 ≻ 626.71 ≪ 843
<i>Bruttel and Kamecke (2012)</i>	590.68	≈	581	≈	592.44	590.68 ≈ 581 ≈ 592.44
<i>Dal Bó (2005)</i>	394.44	≫	365.35	≪	393	412.2 ≻ 394.44 ≻ 365.35 ≪ 393
<i>Dal Bó and Fréchette (2011)</i>	3847.44	≈	3775.96	≈	3802.38	3872.8 ≈ 3847.44 ≈ 3775.96 ≈ 3802.38
<i>Dal Bó and Fréchette (2015)</i>	5522.08	≫	5280.06	≈	5318.33	5549.01 ≻ 5522.08 ≻ 5280.06 ≈ 5318.33
<i>Dreber et al. (2008)</i>	459.11	≈	462.48	≈	481.64	464.92 ≈ 459.11 ≈ 462.48 ≈ 481.64
<i>Duffy and Ochs (2009)</i>	1128.03	≈	1111.23	≈	1123.74	1128.79 ≈ 1128.03 ≈ 1111.23 ≈ 1123.74
<i>Fréchette and Yuksel (2017)</i>	181.98	≫	161.75	<	176.5	186.7 > 181.98 > 161.75 < 176.5
<i>Fudenberg et al. (2012)</i>	356.73	>	331.44	<	347.07	359.72 > 356.73 > 331.44 < 347.07
<i>Kagel and Schley (2013)</i>	1813.04	≈	1862.88	≫	1781.74	1815.98 ≈ 1813.04 ≈ 1862.88 ≻ 1781.74
<i>Sherstyuk et al. (2013)</i>	926.9	≈	953.91	>	912.67	928.91 ≈ 926.9 ≈ 953.91 > 912.67
Pooled	16820.09	≫	16417.97	≪	17165.74	17001.48 ≻ 16820.09 ≻ 16417.97 ≪ 17165.74
Second halves per session						
<i>Aoyagi and Fréchette (2009)</i>	478.46	>	451.68	≈	480.65	480.97 > 478.46 > 451.68 < 480.65
<i>Blonski et al. (2011)</i>	965.51	≫	914.28	≪	1140.5	1005.4 ≻ 965.51 ≻ 914.28 ≪ 1140.5
<i>Bruttel and Kamecke (2012)</i>	342.17	≈	361.38	≈	348.88	343.97 ≈ 342.17 ≈ 361.38 ≈ 348.88
<i>Dal Bó (2005)</i>	464.42	≈	446.63	≪	474.71	474.92 > 464.42 > 446.63 ≪ 474.71
<i>Dal Bó and Fréchette (2011)</i>	3052.84	≈	3081.92	≈	3092.35	3069.11 ≈ 3052.84 ≈ 3081.92 ≈ 3092.35
<i>Dal Bó and Fréchette (2015)</i>	5571.71	>	5419.19	≈	5478.11	5587.96 > 5571.71 > 5419.19 ≈ 5478.11
<i>Dreber et al. (2008)</i>	287.58	≈	285.34	<	303.79	295.44 ≈ 287.58 ≈ 285.34 < 303.79
<i>Duffy and Ochs (2009)</i>	1628.02	≈	1628.89	≈	1609.73	1630.32 ≈ 1628.02 ≈ 1628.89 ≈ 1609.73
<i>Fréchette and Yuksel (2017)</i>	356.12	≈	332.38	<	344.23	362.94 > 356.12 > 332.38 < 344.23
<i>Fudenberg et al. (2012)</i>	443.13	≈	439.28	≈	444.41	445.07 ≈ 443.13 ≈ 439.28 ≈ 444.41
<i>Kagel and Schley (2013)</i>	1213.26	<	1301.76	≫	1203.05	1215.63 < 1213.26 < 1301.76 ≻ 1203.05
<i>Sherstyuk et al. (2013)</i>	587.45	<	640.1	>	597.28	588.7 < 587.45 < 640.1 > 597.28
Pooled	15523.37	≈	15443.11	≪	16124.95	15682.79 ≈ 15523.37 ≈ 15443.11 ≪ 16124.95

Note: Relation signs are used as defined above (Table 7). “No Switching”, “Random Switching” and “Markov Switching” are as defined in the text, but briefly: “No Switching” assumes that each subject randomly chooses a strategy prior to the first supergame and plays this strategy constantly for the entire half session. “Random Switching” assumes that each subject randomly chooses a strategy prior to each supergame (by i.i.d. draws), and “Markov Switching” allows that these switches follow a Markov process. “All but SG” allows subjects to play either AD, Grim, TFT, AC or WSLS, and “Best w/o SG” picks the best mixture model after eliminating AC or WSLS, or both or none of these.

Table 10: Table 9 by treatments – Pure, mixed, or switching strategies?

(a) First halves per session

(b) Second halves per session

Specification	Best w/o SG			All but SG			
	No Switching	Random Switching	Markov Switching	No Switching	Random Switching	Markov Switching	
# Models evaluated	4	4	4	1	1	1	
# Pars estimated (by treatment)	16	16	82	5	5	30	
# Parameters accounted for	3–5	3–5	12–30	5	5	30	
AF09–34	790.2	>	767.22	≈	797.89	≈	797.89
BOS11–9	45.36	≈	43.78	≈	55.54	≈	82.02
BOS11–14	61.42	≈	59.77	≈	72.93	≈	99.82
BOS11–15	29.87	≈	29.93	≈	43.24	≈	70.23
BOS11–16	106.92	≈	103.63	≈	118.03	≈	139.03
BOS11–17	48.66	<	53.62	≈	57.44	≈	82.08
BOS11–26	162.85	≈	140.64	≈	153.41	≈	182.93
BOS11–27	65.77	>	52.93	≈	65.67	≈	89.62
BOS11–30	30.38	≈	30.39	≈	43.76	≈	70.72
BOS11–31	88.41	>	79.66	≈	103.59	≈	113.56
BK12–28	590.68	≈	581	≈	592.44	≈	601.83
D05–18	153.77	≈	136.89	≈	151.01	≈	182.66
D05–19	238.55	≈	226.34	≈	233.48	≈	268.27
DF11–6	379.95	>	357.02	≈	347.58	≈	359.31
DF11–7	549.14	>	513.5	≈	514.91	≈	528.46
DF11–8	581.17	≈	568.86	≈	552.52	≈	562.3
DF11–22	758.46	≈	738.85	≈	742.48	≈	757.73
DF11–23	715.77	≈	735.67	≈	706.54	≈	718.8
DF11–24	843.2	≈	838.63	≈	830.64	≈	839.63
DF15–4	170.82	≈	158.62	≈	164.65	≈	198.4
DF15–5	761.26	≈	724.76	≈	696.32	≈	717.94
DF15–20	1003.49	>	955.28	≈	958.24	≈	973.01
DF15–21	1334.08	≈	1316.87	≈	1274.95	≈	1296.16
DF15–33	1900.51	≈	1796.11	≈	1810.13	≈	1831.38
DF15–35	330.58	>	307.28	≈	319.5	≈	339.6
DRFN08–10	197.04	≈	195.4	≈	208.99	≈	229.14
DRFN08–11	259.69	≈	264.28	≈	264.24	≈	282.44
DO09–32	1128.03	≈	1111.23	≈	1123.74	≈	1134.58
FY17–25	181.98	≈	161.75	<	176.5	≈	199.13
FRD12–29	356.73	>	331.44	<	347.07	≈	366.32
KS13–12	1813.04	≈	1862.88	≈	1781.74	≈	1790.48
STS13–13	926.9	≈	953.91	>	912.67	≈	922.05
<i>Aoyagi and Fréchette (2009)</i>	790.2	>	767.22	≈	797.89	≈	797.89
<i>Blonski et al. (2011)</i>	674.3	≈	626.71	≈	843	≈	1230.45
<i>Bruttel and Kamecke (2012)</i>	590.68	≈	581	≈	592.44	≈	601.83
<i>Dal Bó (2005)</i>	394.44	≈	365.35	≈	393	≈	472.2
<i>Dal Bó and Fréchette (2011)</i>	3847.44	≈	3775.96	≈	3802.38	≈	3927.78
<i>Dal Bó and Fréchette (2015)</i>	5522.08	≈	5280.06	≈	5318.33	≈	5531.46
<i>Dreber et al. (2008)</i>	459.11	≈	462.48	≈	481.64	≈	532.59
<i>Duffy and Ochs (2009)</i>	1128.03	≈	1111.23	≈	1123.74	≈	1134.58
<i>Fréchette and Yuksel (2017)</i>	181.98	≈	161.75	<	176.5	≈	199.13
<i>Fudenberg et al. (2012)</i>	356.73	>	331.44	<	347.07	≈	366.32
<i>Kagel and Schley (2013)</i>	1813.04	≈	1862.88	≈	1781.74	≈	1790.48
<i>Sherstyuk et al. (2013)</i>	926.9	≈	953.91	>	912.67	≈	922.05
Pooled	16820.09	≈	16417.97	≈	17165.74	≈	18601.02

Specification	Best w/o SG			All but SG			
	No Switching	Random Switching	Markov Switching	No Switching	Random Switching	Markov Switching	
# Models evaluated	4	4	4	1	1	1	
# Pars estimated (by treatment)	16	16	82	5	5	30	
# Parameters accounted for	3–5	3–5	12–30	5	5	30	
AF09–34	1266.62	>	1217.36	≈	1251.36	≈	1270.79
BOS11–9	52.34	≈	53.6	<	61.92	≈	85.57
BOS11–14	82.95	≈	99.69	≈	90.37	≈	117.35
BOS11–15	39.65	≈	39.74	≈	53.06	≈	80.05
BOS11–16	179.08	≈	184.91	≈	186.97	≈	209.56
BOS11–17	191.8	≈	190.05	≈	189.58	≈	209.79
BOS11–26	346.61	≈	334.88	≈	355.78	≈	378.67
BOS11–27	213.5	≈	221.89	≈	224.59	≈	235.7
BOS11–30	110.41	≈	100.79	≈	112.27	≈	137.14
BOS11–31	182.25	≈	186.88	≈	197.57	≈	212.21
BK12–28	935.25	≈	944.39	≈	913.55	≈	921.25
D05–18	364.67	≈	363.43	≈	366.23	≈	389.12
D05–19	440.24	≈	433.06	≈	437.46	≈	468.87
DF11–6	672.26	≈	619.89	≈	594.3	≈	594.88
DF11–7	1162.27	≈	1211.18	≈	1107.1	≈	1107.1
DF11–8	1037.38	≈	1024.9	>	942.13	≈	945.91
DF11–22	1465.66	≈	1451.97	≈	1393.15	≈	1409.12
DF11–23	1242.36	≈	1255.79	≈	1201.41	≈	1206.34
DF11–24	1266.48	≈	1253.11	≈	1226.17	≈	1233.69
DF15–4	299.85	>	270.13	≈	272.31	≈	305.5
DF15–5	1390.3	≈	1349.58	>	1265.55	≈	1283.51
DF15–20	1786.19	≈	1792.19	>	1724.15	≈	1724.15
DF15–21	2403.79	≈	2403.31	≈	2310.93	≈	2310.93
DF15–33	4039.13	>	3901.56	>	3804.27	≈	3818.02
DF15–35	978.44	≈	952.81	≈	958.36	≈	976.03
DRFN08–10	308.43	≈	306.9	≈	313.42	≈	330.07
DRFN08–11	436.15	≈	432.7	≈	419.05	≈	436.69
DO09–32	2784.68	≈	2745.23	≈	2702.02	≈	2702.02
FY17–25	504.26	≈	486.46	≈	496.45	≈	514.21
FRD12–29	748.15	≈	750.98	≈	745.6	≈	763.71
KS13–12	3057.93	≈	3175.78	≈	2929.77	≈	2932.51
STS13–13	1500.42	<	1590.96	≈	1457.03	≈	1462.8
<i>Aoyagi and Fréchette (2009)</i>	1266.62	>	1217.36	≈	1251.36	≈	1270.79
<i>Blonski et al. (2011)</i>	1431.75	≈	1426.73	≈	1610.7	≈	1966.5
<i>Bruttel and Kamecke (2012)</i>	935.25	≈	944.39	≈	913.55	≈	921.25
<i>Dal Bó (2005)</i>	807.48	≈	799.06	≈	812.2	≈	879.27
<i>Dal Bó and Fréchette (2011)</i>	6868.86	≈	6841.08	≈	6580.33	≈	6658.59
<i>Dal Bó and Fréchette (2015)</i>	10921.32	>	10694.5	≈	10457.73	≈	10593.11
<i>Dreber et al. (2008)</i>	746.97	≈	742.4	≈	743.19	≈	787.77
<i>Duffy and Ochs (2009)</i>	2784.68	≈	2745.23	≈	2702.02	≈	2702.02
<i>Fréchette and Yuksel (2017)</i>	504.26	≈	486.46	≈	496.45	≈	514.21
<i>Fudenberg et al. (2012)</i>	748.15	≈	750.98	≈	745.6	≈	763.71
<i>Kagel and Schley (2013)</i>	3057.93	≈	3175.78	≈	2929.77	≈	2932.51
<i>Sherstyuk et al. (2013)</i>	1500.42	<	1590.96	≈	1457.03	≈	1462.8
Pooled	31710.95	≈	31564.59	>	31347.14	≈	32546.75

Note: Notation of treatments and meaning of relation signs are all as defined above, see Table 7.

Table 11: Pure, mixed, or switching strategies? Best mixtures without Semi-Grim (ICL-BIC of the models, less is better and relation signs point toward better models)

	Best mixture of pure strategies			Best mixture of generalized pure strategies						
	No Switching	Random Switching	Markov Switching	No Switching	Random Switching	Markov Switching				
Specification										
# Models evaluated	4 ³²	4 ³²	4 ³²	4 ³²	4 ³²	4 ³²				
# Pars estimated (by treatment) (by treatment)	16	16	82	32	32	98				
# Parameters accounted for (by treatment)	3–5	3–5	12–30	6–10	6–10	15–35				
First halves per session										
<i>Aoyagi and Frechette (2009)</i>	790.2	>	767.22	≈	797.89	645.31	≈	646.53	≈	649.53
<i>Blonski et al. (2011)</i>	674.3	≫	626.71	≪	843	713.8	≫	670.62	≪	875.74
<i>Bruttel and Kamecke (2012)</i>	590.68	≈	581	≈	592.44	585.42	≈	570.56	≈	570.87
<i>Dal Bó (2005)</i>	394.44	≫	365.35	≪	393	407.86	≫	378.95	<	404.78
<i>Dal Bó and Fréchette (2011)</i>	3847.44	≈	3775.96	≈	3802.38	3536.73	≈	3589.36	≈	3524.77
<i>Dal Bó and Fréchette (2015)</i>	5522.08	≫	5280.06	≈	5318.33	5259.64	≫	5037.82	≈	5057.54
<i>Dreber et al. (2008)</i>	459.11	≈	462.48	≈	481.64	478.09	≈	466.13	≈	482.86
<i>Duffy and Ochs (2009)</i>	1128.03	≈	1111.23	≈	1123.74	1047.59	≈	1053.04	≈	1049.79
<i>Fréchette and Yuksel (2017)</i>	181.98	≫	161.75	<	176.5	188.5	≈	183.48	≈	175.59
<i>Fudenberg et al. (2012)</i>	356.73	>	331.44	<	347.07	319.45	≈	308.6	≈	320.55
<i>Kagel and Schley (2013)</i>	1813.04	≈	1862.88	≫	1781.74	1761.98	≈	1780.97	>	1694.94
<i>Sherstyuk et al. (2013)</i>	926.9	≈	953.91	>	912.67	865.67	≈	907.14	>	858.65
Pooled	16820.09	≫	16417.97	≪	17165.74	16077.95	>	15853.82	≪	16335.33
Second halves per session										
<i>Aoyagi and Frechette (2009)</i>	478.46	>	451.68	≈	480.65	363.58	≈	368.23	≈	368.89
<i>Blonski et al. (2011)</i>	965.51	≫	914.28	≪	1140.5	992.44	≈	993.71	≪	1154.56
<i>Bruttel and Kamecke (2012)</i>	342.17	≈	361.38	≈	348.88	344.88	≈	358.12	≈	347.08
<i>Dal Bó (2005)</i>	464.42	≈	446.63	≪	474.71	475.11	≈	456.4	≈	469.98
<i>Dal Bó and Fréchette (2011)</i>	3052.84	≈	3081.92	≈	3092.35	2737.11	<	2875.64	>	2721.88
<i>Dal Bó and Fréchette (2015)</i>	5571.71	>	5419.19	≈	5478.11	5164.78	≈	5105.6	≈	5116.42
<i>Dreber et al. (2008)</i>	287.58	≈	285.34	<	303.79	295.06	≈	297.88	≈	303.03
<i>Duffy and Ochs (2009)</i>	1628.02	≈	1628.89	≈	1609.73	1381.01	≈	1416.71	≈	1392.49
<i>Fréchette and Yuksel (2017)</i>	356.12	≈	332.38	<	344.23	309.63	≈	304.7	≈	308.78
<i>Fudenberg et al. (2012)</i>	443.13	≈	439.28	≈	444.41	373.44	≈	395.32	≈	376.62
<i>Kagel and Schley (2013)</i>	1213.26	<	1301.76	≫	1203.05	1170.12	≈	1224.37	>	1143.67
<i>Sherstyuk et al. (2013)</i>	587.45	<	640.1	>	597.28	527.09	≈	590.16	≈	567.63
Pooled	15523.37	≈	15443.11	≪	16124.95	14387.48	<	14656.93	≈	14961.61

Note: Relation signs are used as defined above (Table 7). “No Switching”, “Random Switching” and “Markov Switching” are as defined in the text, but briefly: “No Switching” assumes that each subject randomly chooses a strategy prior to the first supergame and plays this strategy constantly for the entire half session. “Random Switching” assumes that each subject randomly chooses a strategy prior to each supergame (by i.i.d. draws), and “Markov Switching” allows that these switches follow a Markov process. “Best mixture of pure strategies” starts with the general mixture model allowing subjects to play AD, Grim, TFT, AC or WSLS and picks the best-fitting model after eliminating AC or WSLS, or both or none of these. The “Best mixture of generalized strategies” additionally allows for randomization based on these proto-typical strategies as defined in the main text.

Table 13: Best mixtures of pure strategies in relation to a Semi-Grim behavior strategy (ICL-BIC of the models, less is better and relation signs point toward better models)

	Best mixture of pure strategies					Semi-Grim	AD + SG2			
	No Switching	Random Switching	Markov Switching	Best Switching						
Specification										
# Models evaluated	4	4	4			1	1			
# Pars estimated (by treatment)	16	16	82			3	3			
# Parameters accounted for	3–5	3–5	12–35			3	3			
First halves per session										
<i>Aoyagi and Frechette (2009)</i>	790.2	>	767.22	≈	797.89	767.22	≫	694.72	≈	739.24
<i>Blonski et al. (2011)</i>	674.3	≫	626.71	≪	843	626.71	≫	549.45	≪	621.23
<i>Bruttel and Kamecke (2012)</i>	590.68	≈	581	≈	592.44	581	≈	567.86	≈	566
<i>Dal Bó (2005)</i>	394.44	≫	365.35	≪	393	365.35	≈	358.51	<	374.94
<i>Dal Bó and Fréchette (2011)</i>	3847.44	≈	3775.96	≈	3802.38	3775.96	≫	3533.99	≈	3550.1
<i>Dal Bó and Fréchette (2015)</i>	5522.08	≫	5280.06	≈	5318.33	5280.06	≫	4991.74	≈	5014.35
<i>Dreber et al. (2008)</i>	459.11	≈	462.48	≈	481.64	462.48	>	437.17	≈	449.12
<i>Duffy and Ochs (2009)</i>	1128.03	≈	1111.23	≈	1123.74	1111.23	≈	1090.22	≈	1086.17
<i>Fréchette and Yuksel (2017)</i>	181.98	≫	161.75	<	176.5	161.75	≈	161.45	≈	162.3
<i>Fudenberg et al. (2012)</i>	356.73	>	331.44	<	347.07	331.44	>	291.43	≈	295.84
<i>Kagel and Schley (2013)</i>	1813.04	≈	1862.88	≫	1781.74	1862.88	>	1782.82	>	1706.03
<i>Sherstyuk et al. (2013)</i>	926.9	≈	953.91	>	912.67	953.91	≈	912.8	≈	899.24
Pooled	16820.09	≫	16417.97	≪	17165.74	16417.97	≫	15481.59	≈	15573.98
Second halves per session										
<i>Aoyagi and Frechette (2009)</i>	478.46	>	451.68	≈	480.65	478.46	≫	389.24	<	429.54
<i>Blonski et al. (2011)</i>	965.51	≫	914.28	≪	1140.5	965.51	>	867.87	≈	892.86
<i>Bruttel and Kamecke (2012)</i>	342.17	≈	361.38	≈	348.88	342.17	≈	347.4	≈	345.74
<i>Dal Bó (2005)</i>	464.42	≈	446.63	≪	474.71	464.42	>	424.44	≈	430.13
<i>Dal Bó and Fréchette (2011)</i>	3052.84	≈	3081.92	≈	3092.35	3052.84	>	2817.31	≈	2764.96
<i>Dal Bó and Fréchette (2015)</i>	5571.71	>	5419.19	≈	5478.11	5571.71	≫	5043.81	≈	5118.03
<i>Dreber et al. (2008)</i>	287.58	≈	285.34	<	303.79	287.58	≈	264.94	<	288.98
<i>Duffy and Ochs (2009)</i>	1628.02	≈	1628.89	≈	1609.73	1628.02	≫	1403.03	≈	1444.97
<i>Fréchette and Yuksel (2017)</i>	356.12	≈	332.38	<	344.23	356.12	>	313.5	≈	319.17
<i>Fudenberg et al. (2012)</i>	443.13	≈	439.28	≈	444.41	443.13	>	380.75	≈	370.44
<i>Kagel and Schley (2013)</i>	1213.26	<	1301.76	≫	1203.05	1213.26	≈	1211.37	>	1139.73
<i>Sherstyuk et al. (2013)</i>	587.45	<	640.1	>	597.28	587.45	≈	586.72	≈	562.71
Pooled	15523.37	≈	15443.11	≪	16124.95	15523.37	≫	14159.8	≈	14216.67

Note: This table extends Table 9 by picking the best switching model per half-session, after picking the best-fitting mixture involving the pure forms of AD, Grim, TFT, AC and WSLS (as above) for each treatment independently, and examining its goodness-of-fit in relation to Semi-Grim and mixtures involving Semi-Grim.

Table 15: Best mixtures of generalized strategies in relation to a Semi-Grim strategy
(ICL-BIC of the models, less is better and relation signs point toward better models)

	Best mixture of generalized strategies										
	No Switching	Random Switching	Markov Switching	Best Switching	Semi-Grim	AD + SG2					
Specification											
# Models evaluated	4	4	4		1	1					
# Pars estimated (by treatment)	32	32	98		3	3					
# Parameters accounted for	6–10	6–10	15-35		3	3					
First halves per session											
<i>Aoyagi and Frechette (2009)</i>	645.31	≈	646.53	≈	649.53	≈	694.72	≈	739.24		
<i>Blonski et al. (2011)</i>	713.8	≫	670.62	≪	875.74	≈	670.62	≫	549.45	≪	621.23
<i>Bruttel and Kamecke (2012)</i>	585.42	≈	570.56	≈	570.87	≈	570.56	≈	567.86	≈	566
<i>Dal Bó (2005)</i>	407.86	≫	378.95	<	404.78	≈	378.95	>	358.51	<	374.94
<i>Dal Bó and Fréchette (2011)</i>	3536.73	≈	3589.36	≈	3524.77	≈	3589.36	≈	3533.99	≈	3550.1
<i>Dal Bó and Fréchette (2015)</i>	5259.64	≫	5037.82	≈	5057.54	≈	5037.82	≈	4991.74	≈	5014.35
<i>Dreber et al. (2008)</i>	478.09	≈	466.13	≈	482.86	≈	466.13	≫	437.17	≈	449.12
<i>Duffy and Ochs (2009)</i>	1047.59	≈	1053.04	≈	1049.79	≈	1053.04	≈	1090.22	≈	1086.17
<i>Fréchette and Yuksel (2017)</i>	188.5	≈	183.48	≈	175.59	≈	183.48	≫	161.45	≈	162.3
<i>Fudenberg et al. (2012)</i>	319.45	≈	308.6	≈	320.55	≈	308.6	≈	291.43	≈	295.84
<i>Kagel and Schley (2013)</i>	1761.98	≈	1780.97	>	1694.94	≈	1780.97	≈	1782.82	>	1706.03
<i>Sherstyuk et al. (2013)</i>	865.67	≈	907.14	>	858.65	≈	907.14	≈	912.8	≈	899.24
Pooled	16077.95	>	15853.82	≪	16335.33	≈	15853.82	≫	15481.59	≈	15573.98
Second halves per session											
<i>Aoyagi and Frechette (2009)</i>	363.58	≈	368.23	≈	368.89	≈	363.58	≈	389.24	<	429.54
<i>Blonski et al. (2011)</i>	992.44	≈	993.71	≪	1154.56	≈	992.44	≫	867.87	≈	892.86
<i>Bruttel and Kamecke (2012)</i>	344.88	≈	358.12	≈	347.08	≈	344.88	≈	347.4	≈	345.74
<i>Dal Bó (2005)</i>	475.11	≈	456.4	≈	469.98	≈	475.11	≫	424.44	≈	430.13
<i>Dal Bó and Fréchette (2011)</i>	2737.11	<	2875.64	>	2721.88	≈	2737.11	≈	2817.31	≈	2764.96
<i>Dal Bó and Fréchette (2015)</i>	5164.78	≈	5105.6	≈	5116.42	≈	5164.78	≈	5043.81	≈	5118.03
<i>Dreber et al. (2008)</i>	295.06	≈	297.88	≈	303.03	≈	295.06	≈	264.94	<	288.98
<i>Duffy and Ochs (2009)</i>	1381.01	≈	1416.71	≈	1392.49	≈	1381.01	≈	1403.03	≈	1444.97
<i>Fréchette and Yuksel (2017)</i>	309.63	≈	304.7	≈	308.78	≈	309.63	≈	313.5	≈	319.17
<i>Fudenberg et al. (2012)</i>	373.44	≈	395.32	≈	376.62	≈	373.44	≈	380.75	≈	370.44
<i>Kagel and Schley (2013)</i>	1170.12	≈	1224.37	>	1143.67	≈	1170.12	≈	1211.37	>	1139.73
<i>Sherstyuk et al. (2013)</i>	527.09	≈	590.16	≈	567.63	≈	527.09	≈	586.72	≈	562.71
Pooled	14387.48	<	14656.93	≈	14961.61	≈	14387.48	≈	14159.8	≈	14216.67

Note: This table extends Table 11 by picking the best switching model per half-session, after picking the best-fitting mixture involving the generalized forms of AD, Grim, TFT, AC and WLS (as above) for each treatment independently, and examining its goodness-of-fit in relation to Semi-Grim and mixtures involving Semi-Grim.

Table 16: Table 15 by treatments – Best mixtures of generalized strategies in relation to a Semi-Grim strategy

(a) First halves per session

(b) Second halves per session

Specification	Best mixture of generalized strategies				Semi-Grim	AD + SG2
	No Switching	Random Switching	Markov Switching	Best Switching		
# Models evaluated	4	4	4		1	1
# Pars estimated (by treatment)	32	32	98		3	3
# Parameters accounted for	6–10	6–10	15-35		3	3
AF09–34	645.31	≈ 646.53	≈ 649.53	646.53	≈ 694.72	≈ 739.24
BOS11–9	39.68	≈ 34.86	≈ 47.44	34.86	≈ 28.34	≈ 29.99
BOS11–14	52.96	≈ 51.25	≈ 69.19	51.25	≈ 43.74	≈ 56.52
BOS11–15	31.14	≈ 22.09	≈ 45.44	22.09	≈ 14.27	≈ 27.48
BOS11–16	106.93	≈ 109.67	≈ 112.68	109.67	≈ 102.31	≈ 102.4
BOS11–17	56.31	≈ 57.46	≈ 60.52	57.46	≈ 48.33	≈ 53.95
BOS11–26	162.52	≈ 139.66	≈ 158.95	139.66	≈ 130.46	≈ 146.7
BOS11–27	64.41	≈ 54.54	≈ 70.16	54.54	≈ 49.15	≈ 57.95
BOS11–30	42.87	≈ 42.83	≈ 56.32	42.83	≈ 15.97	≈ 30.27
BOS11–31	88.38	≈ 89.66	≈ 94.45	89.66	≈ 86.84	≈ 85.93
BK12–28	585.42	≈ 570.56	≈ 570.87	570.56	≈ 567.86	≈ 566
D05–18	154.75	≈ 138.83	≈ 151.81	138.83	≈ 136.21	≈ 141.56
D05–19	248.85	≈ 235.87	≈ 242.34	235.87	≈ 220.17	≈ 231.25
DF11–6	339.47	≈ 347.55	≈ 316.5	347.55	≈ 333.68	≈ 341.04
DF11–7	526.77	≈ 523.18	≈ 492.09	523.18	≈ 512.99	≈ 525.32
DF11–8	563.75	≈ 552.71	≈ 523.57	552.71	≈ 551.81	≈ 549.33
DF11–22	707.53	≈ 697.75	≈ 694.38	697.75	≈ 684.04	≈ 689.12
DF11–23	619.29	≈ 667.99	≈ 656.5	667.99	≈ 656.29	≈ 643.12
DF11–24	740.3	≈ 758.9	≈ 720.58	758.9	≈ 779.03	≈ 786.01
DF15–4	175.94	≈ 180.54	≈ 168.75	180.54	≈ 149.71	≈ 174.17
DF15–5	737.26	≈ 726.27	≈ 674.63	726.27	≈ 716.51	≈ 689.89
DF15–20	960.95	≈ 913.96	≈ 900.5	913.96	≈ 914.2	≈ 931.26
DF15–21	1234.63	≈ 1210.32	≈ 1218.06	1210.32	≈ 1196.29	≈ 1181.18
DF15–33	1800.02	≈ 1671.27	≈ 1691.77	1671.27	≈ 1714.32	≈ 1735.23
DF15–35	310.13	≈ 293.31	≈ 295.78	293.31	≈ 283.21	≈ 285.12
DRFN08–10	208.83	≈ 197.36	≈ 213.97	197.36	≈ 180.96	≈ 190.57
DRFN08–11	264.48	≈ 263.16	≈ 258.39	263.16	≈ 254.11	≈ 256.45
DO09–32	1047.59	≈ 1053.04	≈ 1049.79	1053.04	≈ 1090.22	≈ 1086.17
FY17–25	188.5	≈ 183.48	≈ 175.59	183.48	≈ 161.45	≈ 162.3
FRD12–29	319.45	≈ 308.6	≈ 320.55	308.6	≈ 291.43	≈ 295.84
KS13–12	1761.98	≈ 1780.97	≈ 1694.94	1780.97	≈ 1782.82	≈ 1706.03
STS13–13	865.67	≈ 907.14	≈ 858.65	907.14	≈ 912.8	≈ 899.24
Aoyagi and Fréchette (2009)	645.31	≈ 646.53	≈ 649.53	646.53	≈ 694.72	≈ 739.24
Blonski et al. (2011)	713.8	≈ 670.62	≈ 875.74	670.62	≈ 549.45	≈ 621.23
Bruttel and Kamecke (2012)	585.42	≈ 570.56	≈ 570.87	570.56	≈ 567.86	≈ 566
Dal Bó (2005)	407.86	≈ 378.95	≈ 404.78	378.95	≈ 358.51	≈ 374.94
Dal Bó and Fréchette (2011)	3536.73	≈ 3589.36	≈ 3524.77	3589.36	≈ 3533.99	≈ 3550.1
Dal Bó and Fréchette (2015)	5259.64	≈ 5037.82	≈ 5057.54	5037.82	≈ 4991.74	≈ 5014.35
Dreber et al. (2008)	478.09	≈ 466.13	≈ 482.86	466.13	≈ 437.17	≈ 449.12
Duffy and Ochs (2009)	1047.59	≈ 1053.04	≈ 1049.79	1053.04	≈ 1090.22	≈ 1086.17
Fréchette and Yuksel (2017)	188.5	≈ 183.48	≈ 175.59	183.48	≈ 161.45	≈ 162.3
Fudenberg et al. (2012)	319.45	≈ 308.6	≈ 320.55	308.6	≈ 291.43	≈ 295.84
Kagel and Schley (2013)	1761.98	≈ 1780.97	≈ 1694.94	1780.97	≈ 1782.82	≈ 1706.03
Sherstyuk et al. (2013)	865.67	≈ 907.14	≈ 858.65	907.14	≈ 912.8	≈ 899.24
Pooled	16077.95	≈ 15853.82	≈ 16335.33	15853.82	≈ 15481.59	≈ 15573.98

Specification	Best mixture of generalized strategies				Semi-Grim	AD + SG2
	No Switching	Random Switching	Markov Switching	Best Switching		
# Models evaluated	4	4	4		1	1
# Pars estimated (by treatment)	32	32	98		3	3
# Parameters accounted for	6–10	6–10	15-35		3	3
AF09–34	998.99	≈ 1011.31	≈ 989.06	989.06	≈ 1085.04	≈ 1178.07
BOS11–9	56.83	≈ 67.33	≈ 63.72	63.72	≈ 50.83	≈ 51.64
BOS11–14	75.9	≈ 70.2	≈ 77.47	77.47	≈ 61.74	≈ 69.39
BOS11–15	43.38	≈ 34.25	≈ 55.25	55.25	≈ 21.88	≈ 37.94
BOS11–16	170.31	≈ 187.28	≈ 176.88	176.88	≈ 178.84	≈ 174.16
BOS11–17	191.55	≈ 192.15	≈ 194.05	194.05	≈ 183.03	≈ 181.55
BOS11–26	328.04	≈ 330.67	≈ 340.16	340.16	≈ 317.33	≈ 317.2
BOS11–27	214.12	≈ 222.33	≈ 218.59	218.59	≈ 231.72	≈ 215.78
BOS11–30	116.18	≈ 109.21	≈ 116.12	116.12	≈ 93.22	≈ 99.81
BOS11–31	183.27	≈ 185.32	≈ 188.11	188.11	≈ 200.87	≈ 188.19
BK12–28	915.48	≈ 926.52	≈ 879.03	879.03	≈ 916.81	≈ 903.59
D05–18	360.73	≈ 361.89	≈ 355.8	355.8	≈ 351.89	≈ 353.01
D05–19	442.98	≈ 435.7	≈ 431.03	431.03	≈ 418.28	≈ 423.88
DF11–6	584.49	≈ 615.02	≈ 531.24	531.24	≈ 599.29	≈ 579.61
DF11–7	1135.74	≈ 1179.37	≈ 996.23	996.23	≈ 1199.27	≈ 1168.14
DF11–8	971.84	≈ 971.32	≈ 874.03	874.03	≈ 985.67	≈ 956.71
DF11–22	1348.86	≈ 1333.85	≈ 1269.49	1269.49	≈ 1331.85	≈ 1303.37
DF11–23	1105.85	≈ 1135.28	≈ 1043.89	1043.89	≈ 1113.94	≈ 1113.42
DF11–24	1076.38	≈ 1105.99	≈ 1038.16	1038.16	≈ 1122.26	≈ 1162.51
DF15–4	298.69	≈ 283.7	≈ 279.23	279.23	≈ 260.02	≈ 283.22
DF15–5	1343.25	≈ 1361.6	≈ 1240.17	1240.17	≈ 1331.49	≈ 1291.93
DF15–20	1696.3	≈ 1689.33	≈ 1589.04	1589.04	≈ 1716.02	≈ 1702.32
DF15–21	2152.39	≈ 2212.29	≈ 2054.9	2054.9	≈ 2213.88	≈ 2176.02
DF15–33	3731.12	≈ 3591.79	≈ 3531.59	3531.59	≈ 3678.61	≈ 3654.5
DF15–35	910.07	≈ 854.66	≈ 870.49	870.49	≈ 856.78	≈ 842.09
DRFN08–10	310.63	≈ 308.06	≈ 311.09	311.09	≈ 270.35	≈ 297.4
DRFN08–11	419.57	≈ 427.47	≈ 404.08	404.08	≈ 425.89	≈ 423.94
DO09–32	2401.04	≈ 2451.43	≈ 2321.8	2321.8	≈ 2526.02	≈ 2569.75
FY17–25	475.52	≈ 454.57	≈ 459.64	459.64	≈ 473.84	≈ 477.72
FRD12–29	636.43	≈ 685.17	≈ 644.13	644.13	≈ 668.18	≈ 652.38
KS13–12	2935.17	≈ 3003.32	≈ 2783.52	2783.52	≈ 3009.28	≈ 2831.74
STS13–13	1358.57	≈ 1493.91	≈ 1325.94	1325.94	≈ 1507.33	≈ 1459.95
Aoyagi and Fréchette (2009)	998.99	≈ 1011.31	≈ 989.06	989.06	≈ 1085.04	≈ 1178.07
Blonski et al. (2011)	1450.51	≈ 1471.96	≈ 1608.54	1608.54	≈ 1369.51	≈ 1365.7
Bruttel and Kamecke (2012)	915.48	≈ 926.52	≈ 879.03	879.03	≈ 916.81	≈ 903.59
Dal Bó (2005)	807.96	≈ 803.27	≈ 797.46	797.46	≈ 772.3	≈ 779.01
Dal Bó and Fréchette (2011)	6260.83	≈ 6383.92	≈ 5874.18	5874.18	≈ 6368.43	≈ 6299.92
Dal Bó and Fréchette (2015)	10177.65	≈ 10037.43	≈ 9693.72	9693.72	≈ 10074.29	≈ 9967.59
Dreber et al. (2008)	734.98	≈ 740.32	≈ 728.28	728.28	≈ 698.35	≈ 723.44
Duffy and Ochs (2009)	2401.04	≈ 2451.43	≈ 2321.8	2321.8	≈ 2526.02	≈ 2569.75
Fréchette and Yuksel (2017)	475.52	≈ 454.57	≈ 459.64	459.64	≈ 473.84	≈ 477.72
Fudenberg et al. (2012)	636.43	≈ 685.17	≈ 644.13	644.13	≈ 668.18	≈ 652.38
Kagel and Schley (2013)	2935.17	≈ 3003.32	≈ 2783.52	2783.52	≈ 3009.28	≈ 2831.74
Sherstyuk et al. (2013)	1358.57	≈ 1493.91	≈ 1325.94	1325.94	≈ 1507.33	≈ 1459.95
Pooled	29421.33	≈ 29743.95	≈ 28851.02	28851.02	≈ 29578.81	≈ 29318.3

Note: Notation of treatments and meaning of relation signs are all as defined above, see Table 7.

Table 17: Best mixtures of pure or generalized strategies in relation to Semi-Grim (ICL-BIC of the models, less is better and relation signs point toward better models)

Specification	Best mixture of pure or generalized strategies								Best Mixture				
	No Switching	Random Switching	Markov Switching	Best Switching	Semi-Grim	SG + Noise	Best Switching	By Treatment					
# Models evaluated	8 ³²	8 ³²	8 ³²	3 × 8 ³²	1	1	24 ³² ≈ 10 ⁴⁴						
# Pars estimated (by treatment)	48	48	180	276	3	4	276						
# Parameters accounted for	3–10	3–10	12–35	3–10	3	4	3–10						
First halves per session													
<i>Aoyagi and Frechette (2009)</i>	645.31	≈	646.53	≈	649.53	≈	646.53	≈	694.72	>	646.73	≈	645.31
<i>Blonski et al. (2011)</i>	677.42	≫	619.43	≪	821.29	≈	619.43	≫	549.45	≈	558.66	<	614.47
<i>Bruttel and Kamecke (2012)</i>	585.42	≈	570.56	≈	570.87	≈	570.56	≈	567.86	≈	559.05	≈	570.56
<i>Dal Bó (2005)</i>	394.44	≫	365.35	≪	393	≈	365.35	≈	358.51	≈	361.22	≈	365.35
<i>Dal Bó and Fréchette (2011)</i>	3536.73	≈	3576.34	≈	3524.77	≈	3576.34	≈	3533.99	>	3416.54	≈	3471.98
<i>Dal Bó and Fréchette (2015)</i>	5250.62	≫	5006.33	≈	5049.55	≈	5006.33	≈	4991.74	≈	4935.93	≈	4965.57
<i>Dreber et al. (2008)</i>	459.11	≈	463.01	≈	477.02	>	463.01	>	437.17	≈	435.65	≈	461.11
<i>Duffy and Ochs (2009)</i>	1047.59	≈	1053.04	≈	1049.79	≈	1053.04	≈	1090.22	>	1017.89	≈	1047.59
<i>Fréchette and Yuksel (2017)</i>	181.98	≫	161.75	<	175.59	≈	161.75	≈	161.45	≈	164.51	≈	161.75
<i>Fudenberg et al. (2012)</i>	319.45	≈	308.6	≈	320.55	≈	308.6	≈	291.43	≈	288.89	≈	308.6
<i>Kagel and Schley (2013)</i>	1761.98	≈	1780.97	>	1694.94	≈	1780.97	≈	1782.82	≈	1726.06	≈	1694.94
<i>Sherstyuk et al. (2013)</i>	865.67	≈	907.14	>	858.65	≈	907.14	≈	912.8	≈	868.89	≈	858.65
Pooled	15951.95	≫	15675.49	≪	16214.1	>	15675.49	>	15481.59	≫	15125.92	≪	15535.6
Second halves per session													
<i>Aoyagi and Frechette (2009)</i>	363.58	≈	368.23	≈	368.89	≈	363.58	≈	389.24	>	353.94	≈	363.58
<i>Blonski et al. (2011)</i>	946.2	>	912.16	≪	1107.23	>	946.2	>	867.87	≈	868.31	≈	908.25
<i>Bruttel and Kamecke (2012)</i>	342.17	≈	358.12	≈	347.08	≈	342.17	≈	347.4	≈	342.6	≈	342.17
<i>Dal Bó (2005)</i>	461.23	≈	445.57	<	469.98	>	461.23	>	424.44	≈	429.06	≈	442.59
<i>Dal Bó and Fréchette (2011)</i>	2737.11	<	2865.45	>	2721.88	≈	2737.11	≈	2817.31	≈	2694.4	≈	2700.09
<i>Dal Bó and Fréchette (2015)</i>	5153.82	≈	5067.01	≈	5106.57	≈	5153.82	≈	5043.81	>	4902.82	≈	4998.09
<i>Dreber et al. (2008)</i>	287.49	≈	281.99	≈	299.36	≈	287.49	≈	264.94	≈	268.86	≈	283.38
<i>Duffy and Ochs (2009)</i>	1381.01	≈	1416.71	≈	1392.49	≈	1381.01	≈	1403.03	>	1308.8	≈	1381.01
<i>Fréchette and Yuksel (2017)</i>	309.63	≈	304.7	≈	308.78	≈	309.63	≈	313.5	≈	278.74	<	304.7
<i>Fudenberg et al. (2012)</i>	373.44	≈	395.32	≈	376.62	≈	373.44	≈	380.75	≈	358.86	≈	373.44
<i>Kagel and Schley (2013)</i>	1170.12	≈	1224.37	>	1143.67	≈	1170.12	≈	1211.37	≈	1153.04	≈	1143.67
<i>Sherstyuk et al. (2013)</i>	527.09	≈	590.16	≈	567.63	≈	527.09	≈	586.72	≈	564.5	≈	527.09
Pooled	14269.01	≈	14448.15	<	14877.81	≈	14269.01	≈	14159.8	≫	13669.82	≪	14108.4

Note: This table extends Table 11 by picking the best switching model per half-session, after picking the best-fitting mixture involving either the pure or generalized forms of AD, Grim, TFT, AC and WSLS (as above) for each treatment independently, and examining its goodness-of-fit in relation to Semi-Grim and mixtures involving Semi-Grim.

Table 18: Table 17 by treatments – Best mixtures of pure or generalized strategies in relation to Semi-Grim

(a) First halves per session

Specification	Best mixture of pure or generalized strategies				Semi-Grim	SG + Noise	Best Mixture Best Switching By Treatment
	No Switching	Random Switching	Markov Switching	Best Switching			
# Models evaluated	8^{32}	8^{32}	8^{32}	3×8^{32}	1	1	$24^{32} \approx 10^{44}$
# Pars estimated (by treatment)	48	48	180	276	3	4	276
# Parameters accounted for	3–10	3–10	12–35	3–10	3	4	3–10
AF09–34	645.31	646.53	649.53	646.53	694.72	646.73	645.31
BOS11–9	39.68	34.86	47.44	34.86	28.34	28.33	34.86
BOS11–14	52.96	51.25	69.19	51.25	43.74	45.24	51.25
BOS11–15	29.87	22.09	43.24	22.09	14.27	15.77	22.09
BOS11–16	106.92	103.63	112.68	103.63	102.31	93.04	103.63
BOS11–17	48.66	53.62	57.44	53.62	48.33	54.96	48.66
BOS11–26	162.52	139.66	153.41	139.66	130.46	128.89	139.66
BOS11–27	64.41	52.93	65.67	52.93	49.15	50.65	52.93
BOS11–30	30.38	30.39	43.76	30.39	15.97	17.47	30.38
BOS11–31	88.38	79.66	94.45	79.66	86.84	84.25	79.66
BK12–28	585.42	570.56	570.87	570.56	567.86	559.05	570.56
D05–18	153.77	136.89	151.01	136.89	136.21	139.6	136.89
D05–19	238.55	226.34	233.48	226.34	220.17	218.78	226.34
DF11–6	339.47	347.55	347.55	347.55	333.68	318.37	316.5
DF11–7	526.77	513.5	492.09	513.5	512.99	501.16	492.09
DF11–8	563.75	552.71	523.57	552.71	551.81	547.22	523.57
DF11–22	707.53	697.75	694.38	697.75	684.04	655.17	694.38
DF11–23	619.29	667.99	656.5	667.99	656.29	639.07	619.29
DF11–24	740.3	758.9	720.58	758.9	779.03	734	720.58
DF15–4	170.82	158.62	164.65	158.62	149.71	151.68	158.62
DF15–5	737.26	724.76	674.63	724.76	716.51	705.07	674.63
DF15–20	960.95	913.96	900.5	913.96	914.2	901.08	900.5
DF15–21	1234.63	1210.32	1218.06	1210.32	1196.29	1174.44	1210.32
DF15–33	1800.02	1671.27	1691.77	1671.27	1714.32	1703.43	1671.27
DF15–35	310.13	293.31	295.78	293.31	283.21	276.9	293.31
DRFN08–10	197.04	195.4	208.99	195.4	180.96	182.62	195.4
DRFN08–11	259.69	263.16	258.39	263.16	254.11	250.23	258.39
DO09–32	1047.59	1053.04	1049.79	1053.04	1090.22	1017.89	1047.59
FY17–25	181.98	161.75	175.59	161.75	161.45	164.51	161.75
FRD12–29	319.45	308.6	320.55	308.6	291.43	288.89	308.6
KS13–12	1761.98	1780.97	1694.94	1780.97	1782.82	1726.06	1694.94
STS13–13	865.67	907.14	858.65	907.14	912.8	868.89	858.65
<i>Aoyagi and Fréchette (2009)</i>	645.31	646.53	649.53	646.53	694.72	646.73	645.31
<i>Blonski et al. (2011)</i>	677.42	619.43	821.29	619.43	549.45	558.66	614.47
<i>Brutel and Kamecke (2012)</i>	585.42	570.56	570.87	570.56	567.86	559.05	570.56
<i>Dal Bó (2005)</i>	394.44	365.35	393	365.35	358.51	361.22	365.35
<i>Dal Bó and Fréchette (2011)</i>	3536.73	3576.34	3524.77	3576.34	3533.99	3416.54	3471.98
<i>Dal Bó and Fréchette (2015)</i>	5250.62	5006.33	5049.55	5006.33	4991.74	4935.93	4965.57
<i>Dreber et al. (2008)</i>	459.11	463.01	477.02	463.01	437.17	435.65	461.11
<i>Duffy and Ochs (2009)</i>	1047.59	1053.04	1049.79	1053.04	1090.22	1017.89	1047.59
<i>Fréchette and Yuksel (2017)</i>	181.98	161.75	175.59	161.75	161.45	164.51	161.75
<i>Fudenberg et al. (2012)</i>	319.45	308.6	320.55	308.6	291.43	288.89	308.6
<i>Kagel and Schley (2013)</i>	1761.98	1780.97	1694.94	1780.97	1782.82	1726.06	1694.94
<i>Sherstyuk et al. (2013)</i>	865.67	907.14	858.65	907.14	912.8	868.89	858.65
Pooled	15951.95	15675.49	16214.1	15675.49	15481.59	15125.92	15535.6

(b) Second halves per session

Specification	Best mixture of pure or generalized strategies				Semi-Grim	SG + Noise	Best Mixture Best Switching By Treatment
	No Switching	Random Switching	Markov Switching	Best Switching			
# Models evaluated	8^{32}	8^{32}	8^{32}	3×8^{32}	1	1	$24^{32} \approx 10^{44}$
# Pars estimated (by treatment)	48	48	180	276	3	4	276
# Parameters accounted for	3–10	3–10	12–35	3–10	3	4	3–10
AF09–34	998.99	1011.31	989.06	989.06	1085.04	1002.44	989.06
BOS11–9	52.34	53.6	61.92	53.6	50.83	52.33	52.34
BOS11–14	75.9	70.2	77.47	77.47	61.78	63.23	70.2
BOS11–15	39.65	34.25	53.06	53.06	21.88	24.13	34.25
BOS11–16	170.31	184.91	176.88	176.88	178.84	159.07	170.31
BOS11–17	191.55	190.05	189.58	189.58	183.03	184.58	189.58
BOS11–26	328.04	330.67	340.16	340.16	317.33	315.97	328.04
BOS11–27	213.5	221.89	218.59	218.59	231.72	229.23	213.5
BOS11–30	110.41	100.79	112.27	112.27	93.22	94.72	100.79
BOS11–31	182.25	185.32	188.11	188.11	200.87	188.4	182.25
BK12–28	915.48	926.52	879.03	879.03	916.81	902.34	879.03
D05–18	360.73	361.89	355.8	355.8	351.89	352.95	355.8
D05–19	440.24	433.06	431.03	431.03	418.28	409	431.03
DF11–6	584.49	615.02	531.24	531.24	599.29	543.48	531.24
DF11–7	1135.74	1179.37	996.23	996.23	1199.27	1089.22	996.23
DF11–8	971.84	971.32	874.03	874.03	985.67	985.55	874.03
DF11–22	1348.86	1333.85	1269.49	1269.49	1331.85	1333.74	1269.49
DF11–23	1105.85	1135.28	1043.89	1043.89	1113.94	1098.69	1043.89
DF11–24	1076.38	1105.99	1038.16	1038.16	1122.26	1088.71	1038.16
DF15–4	298.69	270.13	272.31	272.31	260.02	262.05	270.13
DF15–5	1343.25	1349.58	1240.17	1240.17	1331.49	1325.96	1240.17
DF15–20	1696.3	1689.33	1589.04	1589.04	1716.02	1677.48	1589.04
DF15–21	2152.39	2212.29	2054.9	2054.9	2213.88	2121.22	2054.9
DF15–33	3731.12	3591.79	3531.59	3531.59	3678.61	3633.99	3531.59
DF15–35	910.07	854.66	870.49	870.49	856.78	858.57	854.66
DRFN08–10	308.43	306.9	311.09	311.09	270.35	272.02	306.9
DRFN08–11	419.57	427.47	404.08	404.08	425.89	427.44	404.08
DO09–32	2401.04	2451.43	2321.8	2321.8	2526.02	2344.75	2321.8
FY17–25	475.52	454.57	459.64	459.64	473.84	451.38	454.57
FRD12–29	636.43	685.17	644.13	644.13	668.18	632.09	636.43
KS13–12	2935.17	3003.32	2783.52	2783.52	3009.28	2904.04	2783.52
STS13–13	1358.57	1493.91	1325.94	1325.94	1507.33	1446.27	1325.94
<i>Aoyagi and Fréchette (2009)</i>	998.99	1011.31	989.06	989.06	1085.04	1002.44	989.06
<i>Blonski et al. (2011)</i>	1413	1425.33	1572.06	1572.06	1369.51	1351.72	1402.96
<i>Brutel and Kamecke (2012)</i>	915.48	926.52	879.03	879.03	916.81	902.34	879.03
<i>Dal Bó (2005)</i>	804.43	799.3	797.46	797.46	772.3	764.79	797.46
<i>Dal Bó and Fréchette (2011)</i>	6260.83	6383.92	5874.18	5874.18	6368.43	6160.93	5874.18
<i>Dal Bó and Fréchette (2015)</i>	10177.65	10006.1	9682.9	9682.9	10074.29	9902.6	9631.07
<i>Dreber et al. (2008)</i>	731.63	737.99	728.28	728.28	698.35	702.26	718.29
<i>Duffy and Ochs (2009)</i>	2401.04	2451.43	2321.8	2321.8	2526.02	2344.75	2321.8
<i>Fréchette and Yuksel (2017)</i>	475.52	454.57	459.64	459.64	473.84	451.38	454.57
<i>Fudenberg et al. (2012)</i>	636.43	685.17	644.13	644.13	668.18	632.09	636.43
<i>Kagel and Schley (2013)</i>	2935.17	3003.32	2783.52	2783.52	3009.28	2904.04	2783.52
<i>Sherstyuk et al. (2013)</i>	1358.57	1493.91	1325.94	1325.94	1507.33	1446.27	1325.94
Pooled	29345.07	29624.32	28779.62	28779.62	29578.81	28711.49	28340.25

Note: Notation of treatments and meaning of relation signs are all as defined above, see Table 7.

Table 19: Is there a single “semi grim” type? Mixture models involving SG

	SG + TFT		SG + AD		3 × SG		2 × SG + Noise		SG + SG		SG + Noise		Semi-Grim
Specification													
# Models evaluated	1		1		1		1		1		1		1
# Pars estimated (by treatment)	5		5		11		8		7		4		3
# Parameters accounted for	5		5		11		8		7		4		3
First halves per session													
<i>Aoyagi and Frechette (2009)</i>	698.37	≈	698.36	>	650.89	≈	647.47	≈	650.83	≈	646.73	≈	694.72
<i>Blonski et al. (2011)</i>	614.21	≪	663.28	≪	743.15	≫	667.95	≫	642.5	≫	558.66	≈	549.45
<i>Bruttel and Kamecke (2012)</i>	571.1	≈	567.84	≈	556.46	≈	554.86	≈	553.59	≈	559.05	≈	567.86
<i>Dal Bó (2005)</i>	386.98	≈	378.87	≈	401.23	≫	385.12	≈	379.63	>	361.22	≈	358.51
<i>Dal Bó and Fréchette (2011)</i>	3538.27	≈	3488.81	>	3327.75	≈	3325.99	≈	3335.4	<	3416.54	<	3533.99
<i>Dal Bó and Fréchette (2015)</i>	5047.74	≈	5024.63	≈	5004.06	>	4946.44	≈	4958.73	≈	4935.93	≈	4991.74
<i>Dreber et al. (2008)</i>	443.9	≈	442.47	≈	454.03	≈	442.03	≈	438.35	≈	435.65	≈	437.17
<i>Duffy and Ochs (2009)</i>	1099.57	≈	1084.74	>	1019.33	≈	1018.29	≈	1016.91	≈	1017.89	<	1090.22
<i>Fréchette and Yuksel (2017)</i>	180.54	≫	165.94	≪	187.62	≈	183.63	≈	179.6	>	164.51	≈	161.45
<i>Fudenberg et al. (2012)</i>	297.99	≈	292.47	≈	290.04	≈	289.31	≈	287.31	≈	288.89	≈	291.43
<i>Kagel and Schley (2013)</i>	1777.58	>	1689.17	≈	1622.36	≈	1680.43	≈	1677.99	≈	1726.06	≈	1782.82
<i>Sherstyuk et al. (2013)</i>	904.29	≈	887.79	≈	868.66	≈	870.16	≈	868.11	≈	868.89	≈	912.8
Pooled	15742.91	>	15566.74	≈	15526.79	≫	15303.48	≈	15244.28	≈	15125.92	≪	15481.59
Second halves per session													
<i>Aoyagi and Frechette (2009)</i>	397.94	≈	392.87	≈	363.39	>	360.1	≈	358.14	>	353.94	≈	389.24
<i>Blonski et al. (2011)</i>	940.75	≈	914.98	<	985.46	≫	927.39	≫	906.87	≈	868.31	≈	868.72
<i>Bruttel and Kamecke (2012)</i>	337.22	≈	349.32	≫	312.72	≈	316.23	≈	322.87	≈	342.6	≈	347.4
<i>Dal Bó (2005)</i>	444.73	≈	437.71	≪	476.01	>	458.8	≈	455.9	≫	429.06	≈	424.44
<i>Dal Bó and Fréchette (2011)</i>	2798.46	>	2640.88	≈	2554.16	≈	2556.55	≈	2545.23	≪	2694.4	<	2821.38
<i>Dal Bó and Fréchette (2015)</i>	5053.74	≈	5032.59	≈	4946.96	≈	4909.92	≈	4879.57	≈	4902.82	<	5043.81
<i>Dreber et al. (2008)</i>	262.93	≈	275.84	≈	283.45	>	272.19	≈	268.1	≈	268.86	≈	265.11
<i>Duffy and Ochs (2009)</i>	1414.84	≈	1399.65	>	1288.79	≈	1284.98	≈	1284.59	≈	1308.8	<	1403.03
<i>Fréchette and Yuksel (2017)</i>	319.44	≈	319.78	≈	283.19	≈	278.98	≈	277.03	≈	278.74	<	313.5
<i>Fudenberg et al. (2012)</i>	388.9	>	366.79	≈	345.11	≈	345.11	≈	348.86	≈	358.86	≈	380.75
<i>Kagel and Schley (2013)</i>	1200.87	>	1126.89	≈	1065.1	≈	1070.99	≈	1111.19	≈	1153.04	≈	1211.37
<i>Sherstyuk et al. (2013)</i>	567.48	≈	550.64	≈	518	≈	518.83	≈	544.85	≈	564.5	≈	586.72
Pooled	14309.69	≫	13990.32	≈	13823.56	≫	13591.86	≈	13558.53	≈	13669.82	≪	14164.9

Note: This table verifies a number of possible mixtures involving Semi-Grim types as a robustness check for the sufficiency of focusing on the mixtures examined above. E.g. “3 × SG refers to a model containing three different versions of memory-1 semi-grim with allowing for heterogeneity of randomization parameters across subjects.

Table 20: Table 19 by treatments – Is there a single “semi grim” type? Mixture models involving SG

(a) First halves per session

(b) Second halves per session

	SG + TFT	SG + AD	3× SG	2× SG + Noise	SG + SG	SG + Noise	Semi-Grim
Specification	1	1	1	1	1	1	1
# Models evaluated	5	5	11	8	7	4	3
# Pars estimated (by treatment)	5	5	11	8	7	4	3
# Parameters accounted for	5	5	11	8	7	4	3
AF09-34	698.37	698.36	650.89	647.47	650.83	646.73	694.72
BOS11-9	37.72	36.71	36.42	31.91	30.41	28.33	28.34
BOS11-14	46.74	59.53	62.1	57.42	55.89	45.24	43.74
BOS11-15	17.7	30.49	34.31	29.77	28.27	15.77	14.27
BOS11-16	105.34	108.05	102.21	97.29	95.79	93.04	102.31
BOS11-17	51.6	54.39	62.12	58.67	55.43	54.96	48.33
BOS11-26	138.93	145.92	143.32	135.62	133.77	128.89	130.46
BOS11-27	57.08	57.83	66.2	61.51	60.21	50.65	49.15
BOS11-30	19.19	31.51	42.82	36.78	35.23	17.47	15.97
BOS11-31	89.84	88.76	83.49	78.88	77.39	84.25	86.84
BK12-28	571.1	567.84	556.46	554.86	553.59	559.05	567.86
D05-18	152.64	140.51	147.72	137.32	135.38	139.6	136.21
D05-19	230.8	234.82	245.71	242.13	239.28	218.78	220.17
DF11-6	340.51	342.93	318	315.22	312.8	318.37	333.68
DF11-7	514.24	510.31	497	493.05	489.18	501.16	512.99
DF11-8	555.59	521.04	510.54	513.57	511.26	545.22	551.81
DF11-22	683.89	690.37	638.31	641.64	649.47	677.12	684.04
DF11-23	652.77	625.94	589.27	588.49	605.97	639.07	656.29
DF11-24	764.34	771.29	715.4	730.93	729.03	734	779.03
DF15-4	160.08	174.81	190.8	179.18	177.25	151.68	149.71
DF15-5	729.79	693.23	694.48	688.23	684.87	705.07	716.51
DF15-20	915.67	936	903.1	898.57	895.31	901.08	914.2
DF15-21	1214.46	1167.96	1149.36	1139.86	1173.16	1174.44	1196.29
DF15-33	1713.96	1740.36	1712.55	1708.28	1705.28	1703.43	1714.32
DF15-35	284.62	283.12	289.62	285.66	282.04	276.9	283.21
DRFN08-10	185.08	185.52	190.34	185.08	183.64	182.62	180.96
DRFN08-11	255.32	253.45	255.98	251.35	249.8	250.23	254.11
DO09-32	1099.57	1084.74	1019.33	1018.29	1016.91	1017.89	1090.22
FY17-25	180.54	165.94	187.62	183.63	179.6	164.51	161.45
FRD12-29	297.99	292.47	290.04	289.31	287.31	288.89	291.43
KS13-12	1777.58	1689.17	1622.36	1680.43	1677.99	1726.06	1782.82
STS13-13	904.29	887.79	868.66	870.16	868.11	868.89	912.8
<i>Aoyagi and Fréchet (2009)</i>	698.37	698.36	650.89	647.47	650.83	646.73	694.72
<i>Blonski et al. (2011)</i>	614.21	663.28	743.15	667.95	642.5	558.66	549.45
<i>Brutzel and Kamecke (2012)</i>	571.1	567.84	556.46	554.86	553.59	559.05	567.86
<i>Dal Bó (2005)</i>	386.98	378.87	401.23	385.12	379.63	361.22	358.51
<i>Dal Bó and Fréchet (2011)</i>	3538.27	3488.81	3327.75	3325.99	3335.4	3416.54	3533.99
<i>Dal Bó and Fréchet (2015)</i>	5047.74	5024.63	5004.06	4946.44	4958.73	4935.93	4991.74
<i>Dreber et al. (2008)</i>	443.9	442.47	454.03	442.03	438.35	435.65	437.17
<i>Duffy and Ochs (2009)</i>	1099.57	1084.74	1019.33	1018.29	1016.91	1017.89	1090.22
<i>Fréchet and Yuksel (2017)</i>	180.54	165.94	187.62	183.63	179.6	164.51	161.45
<i>Fudenberg et al. (2012)</i>	297.99	292.47	290.04	289.31	287.31	288.89	291.43
<i>Kagel and Schley (2013)</i>	1777.58	1689.17	1622.36	1680.43	1677.99	1726.06	1782.82
<i>Sherstyuk et al. (2013)</i>	904.29	887.79	868.66	870.16	868.11	868.89	912.8
Pooled	15742.91	15566.74	15526.79	15303.48	15244.28	15125.92	15481.59

	SG + TFT	SG + AD	3× SG	2× SG + Noise	SG + SG	SG + Noise	Semi-Grim
Specification	1	1	1	1	1	1	1
# Models evaluated	5	5	11	8	7	4	3
# Pars estimated (by treatment)	5	5	11	8	7	4	3
# Parameters accounted for	5	5	11	8	7	4	3
AF09-34	1089.07	1088.68	990.69	991.1	1004.14	1002.44	1085.04
BOS11-9	53.89	53.88	64.12	59.61	58.1	52.33	50.83
BOS11-14	64.87	72.35	81.52	76.99	75.39	63.23	61.74
BOS11-15	27.62	40.96	39.5	35	33.53	24.13	22.64
BOS11-16	183.6	177.15	166.77	162.08	160.58	159.07	178.84
BOS11-17	180.21	181.35	190.99	186.36	184.86	184.58	183.03
BOS11-26	316.74	303.4	295.8	303.64	301.8	315.97	317.33
BOS11-27	236.3	216.62	206.64	211.44	211.81	229.23	231.72
BOS11-30	96.22	102.16	106.32	101.82	100.32	94.72	93.22
BOS11-31	205.18	191.19	173.37	168.88	167.41	188.4	200.87
BK12-28	908.74	904.92	865.47	871.79	884.68	902.34	916.81
D05-18	355.5	357.75	371.46	365.36	362.37	352.95	351.89
D05-19	424.87	426.93	429.94	423.33	420.68	409	418.28
DF11-6	592.73	550.5	505.61	531.1	529.38	543.48	599.29
DF11-7	1175.92	1122.24	1076.72	1033.39	1080.6	1089.22	1199.27
DF11-8	998.86	925.57	888.74	883.02	881.18	985.55	985.67
DF11-22	1299.41	1280.15	1252.38	1259.53	1257.63	1333.74	1313.85
DF11-23	1107.01	1076.91	1018.94	1017.54	1030.94	1098.69	1113.94
DF11-24	1098.92	1126.04	1031.22	1031.12	1029.26	1088.71	1122.26
DF15-4	260.03	276.63	278.3	271.52	269.36	262.05	260.02
DF15-5	1343.78	1288.52	1263.24	1261.59	1261.3	1325.96	1331.49
DF15-20	1685.68	1666.64	1638.48	1641.52	1658.8	1677.48	1716.02
DF15-21	2221.3	2133.97	2092.08	2094.42	2097.89	2121.22	2213.88
DF15-33	3602.23	3627.79	3508.8	3550.24	3547.85	3633.99	3678.61
DF15-35	845.37	823.8	841.54	836.16	834.42	858.57	856.78
DRFN08-10	273.7	274.83	269.72	264.73	263.18	272.02	270.35
DRFN08-11	423.75	425.89	404.27	405.55	404	427.44	425.89
DO09-32	2524.24	2510.62	2305.09	2300.39	2311.91	2344.75	2526.02
FY17-25	476.73	476.04	444.24	441.83	443.06	451.38	473.84
FRD12-29	657.64	645.3	589.46	608.61	613.78	632.09	668.18
KS13-12	2984.62	2795.09	2642.29	2711.18	2804.08	2904.04	3009.28
STS13-13	1461.59	1429.53	1368.44	1411.21	1430.92	1446.27	1507.33
<i>Aoyagi and Fréchet (2009)</i>	1089.07	1088.68	990.69	991.1	1004.14	1002.44	1085.04
<i>Blonski et al. (2011)</i>	1414.72	1389.14	1435.21	1385.94	1363.89	1351.72	1370.26
<i>Brutzel and Kamecke (2012)</i>	908.74	904.92	865.47	871.79	884.68	902.34	916.81
<i>Dal Bó (2005)</i>	783.92	788.23	809.2	794.36	788.01	764.79	772.3
<i>Dal Bó and Fréchet (2011)</i>	6299.78	6108.34	5832.86	5798.78	5846.68	6160.93	6368.43
<i>Dal Bó and Fréchet (2015)</i>	9987.54	9846.51	9686.58	9702.12	9710.44	9902.6	10074.29
<i>Dreber et al. (2008)</i>	700.95	704.23	681.69	675.88	672.08	702.26	698.35
<i>Duffy and Ochs (2009)</i>	2524.24	2510.62	2305.09	2300.39	2311.91	2344.75	2526.02
<i>Fréchet and Yuksel (2017)</i>	476.73	476.04	444.24	441.83	443.06	451.38	473.84
<i>Fudenberg et al. (2012)</i>	657.64	645.3	589.46	608.61	613.78	632.09	668.18
<i>Kagel and Schley (2013)</i>	2984.62	2795.09	2642.29	2711.18	2804.08	2904.04	3009.28
<i>Sherstyuk et al. (2013)</i>	1461.59	1429.53	1368.44	1411.21	1430.92	1446.27	1507.33
Pooled	29471.93	28869	28052.44	27984.99	28128.99	28711.49	29579.56

Note: Notation of treatments and meaning of relation signs are all as defined above, see Table 7.

Appendix D: Robustness checks for Section 4

The tables in this section replicate the tables presented in the Section 4, again provide robustness checks and the results treatment-by-treatment.

- Table 21 shows aggregate state-wise cooperation rates for different lagged histories (cooperation or defection of the opponent in $t - 2$) *TFT-Scheme*.
- Table 22 shows aggregate state-wise cooperation rates for different lagged histories (joint cooperation or not in $t - 2$) *Grim-Scheme*.
- Table 24 compares different models containing semi-grim to models containing pure strategies assuming no-switching behavior.
- Table 26 compares different models containing semi-grim to models containing pure strategies assuming random-switching behavior.
- Table 28 compares different models containing modifications of semi-grim.
- Table 30 compares different models containing prototypical strategies derived from strategies discussed in previous literature in a No-Switching model.
- Table 32 compares different two parameter versions of semi-grim with models containing prototypical strategies. The memory-2 level follows a *Grim-Scheme* if applicable
- Table 33 compares different two parameter versions of semi-grim with models containing prototypical strategies. The memory-2 level follows a *TFT-Scheme* if applicable
- Table 34 examines all mixtures of Semi-Grim with pure or generalized pure strategies as secondary components (robustness check for Table 4).

Table 21: Strategies as a function of behavior in $t - 2$ (TFT scheme)

Experiment	Cooperation after $\emptyset, (c, c), (d, c)$ in $t - 2$				Cooperation after $(c, d), (d, d)$ in $t - 2$			
	$\hat{\delta}_{cc}$	$\hat{\delta}_{cd}$	$\hat{\delta}_{dc}$	$\hat{\delta}_{dd}$	$\hat{\delta}_{cc}$	$\hat{\delta}_{cd}$	$\hat{\delta}_{dc}$	$\hat{\delta}_{dd}$
First halves per session								
<i>Aoyagi and Frechette (2009)</i>	0.93	>> 0.439	≈ 0.388	≈ 0.434	0.789	>> 0.463	≈ 0.44	> 0.291
<i>Blonski et al. (2011)</i>	0.901	>> 0.27	≈ 0.146	>> 0.053	0.667	≈ 0.296	≈ 0.321	>> 0.027
<i>Bruttel and Kamecke (2012)</i>	0.908	>> 0.312	≈ 0.218	≈ 0.151	0.944	>> 0.247	≈ 0.247	>> 0.063
<i>Dal Bó (2005)</i>	0.93	>> 0.232	≈ 0.31	> 0.126	0.833	> 0.147	≈ 0.413	>> 0.071
<i>Dal Bó and Fréchet (2011)</i>	0.955	>> 0.352	≈ 0.298	>> 0.086	0.885	>> 0.291	≈ 0.41	>> 0.048
<i>Dal Bó and Fréchet (2015)</i>	0.944	>> 0.301	≈ 0.277	>> 0.098	0.847	>> 0.288	≈ 0.44	>> 0.044
<i>Dreber et al. (2008)</i>	0.902	>> 0.213	≈ 0.189	>> 0.061	1	> 0.233	≈ 0.302	>> 0.025
<i>Duffy and Ochs (2009)</i>	0.927	>> 0.316	≈ 0.304	≈ 0.232	0.691	>> 0.277	≈ 0.361	>> 0.08
<i>Fréchet and Yuksel (2017)</i>	0.943	>> 0.153	≈ 0.241	≈ 0.1	1	≈	≈ 0.4	≈ 0.086
<i>Fudenberg et al. (2012)</i>	0.984	>> 0.394	≈ 0.347	>> 0.05	0.895	>> 0.41	≈ 0.579	>> 0.069
<i>Kagel and Schley (2013)</i>	0.94	>> 0.29	≈ 0.25	>> 0.125	0.787	>> 0.196	≈ 0.402	>> 0.032
<i>Sherstyuk et al. (2013)</i>	0.951	>> 0.329	≈ 0.341	> 0.186	0.844	>> 0.328	≈ 0.424	>> 0.09
Pooled	0.944	>> 0.312	> 0.279	>> 0.106	0.826	>> 0.287	≈ 0.41	>> 0.05
Second halves per session								
<i>Aoyagi and Frechette (2009)</i>	0.961	>> 0.408	≈ 0.567	≈ 0.447	0.867	>> 0.381	≈ 0.451	≈ 0.328
<i>Blonski et al. (2011)</i>	0.922	>> 0.224	≈ 0.195	>> 0.029	0.944	>> 0.402	≈ 0.324	>> 0.018
<i>Bruttel and Kamecke (2012)</i>	0.948	>> 0.239	≈ 0.214	≈ 0.118	0.923	> 0.167	≈ 0.5	>> 0.018
<i>Dal Bó (2005)</i>	0.919	>> 0.264	≈ 0.39	>> 0.113	0.938	>> 0.175	≈ 0.383	>> 0.047
<i>Dal Bó and Fréchet (2011)</i>	0.979	>> 0.391	≈ 0.29	>> 0.075	0.975	>> 0.334	≈ 0.547	>> 0.022
<i>Dal Bó and Fréchet (2015)</i>	0.977	>> 0.304	≈ 0.328	>> 0.064	0.927	>> 0.343	≈ 0.532	>> 0.028
<i>Dreber et al. (2008)</i>	0.917	>> 0.111	< 0.311	>> 0.005	0.909	> 0.5	≈ 0.629	>> 0.01
<i>Duffy and Ochs (2009)</i>	0.98	>> 0.408	≈ 0.371	> 0.232	0.849	>> 0.316	≈ 0.415	>> 0.058
<i>Fréchet and Yuksel (2017)</i>	0.973	>> 0.213	≈ 0.286	≈ 0.214	0.818	≈ 0.286	≈ 0.575	>> 0.038
<i>Fudenberg et al. (2012)</i>	0.974	>> 0.5	≈ 0.41	>> 0.111	0.84	> 0.463	≈ 0.417	>> 0.075
<i>Kagel and Schley (2013)</i>	0.967	>> 0.281	≈ 0.263	>> 0.061	0.872	>> 0.188	≈ 0.527	>> 0.018
<i>Sherstyuk et al. (2013)</i>	0.973	>> 0.503	≈ 0.417	>> 0.12	0.968	>> 0.431	≈ 0.5	>> 0.062
Pooled	0.973	>> 0.325	≈ 0.315	>> 0.076	0.917	>> 0.332	≈ 0.499	>> 0.028

Note: Relation signs, bootstrap procedure, and derived p -values are exactly as above, see Table 1, with the obvious adaptation that the Holm-Bonferroni correction now applies to all eight tests per data set.

Table 22: Strategies as a function of behavior in $t - 2$ (Grim scheme)

Experiment	Cooperation after $\emptyset, (c, c)$ in $t - 2$				Cooperation after $(c, d), (d, c), (d, d)$ in $t - 2$								
	$\hat{\delta}_{cc}$	$\hat{\delta}_{cd}$	$\hat{\delta}_{dc}$	$\hat{\delta}_{dd}$	$\hat{\delta}_{cc}$	$\hat{\delta}_{cd}$	$\hat{\delta}_{dc}$	$\hat{\delta}_{dd}$					
First halves per session													
<i>Aoyagi and Frechette (2009)</i>	0.939	>> 0.39	≈	0.439	≈	0.556	0.782	>> 0.485	≈	0.39	>	0.32	
<i>Blonski et al. (2011)</i>	0.903	>> 0.248	≈	0.174	>>	0.045	0.714	>	0.318	≈	0.216	>	0.031
<i>Bruttel and Kamecke (2012)</i>	0.919	>> 0.296	≈	0.245	≈	0.179	0.833	>>	0.278	≈	0.213	>>	0.071
<i>Dal Bó (2005)</i>	0.926	>> 0.184	≈	0.31	≈	0.143	0.889	>>	0.254	≈	0.39	>>	0.074
<i>Dal Bó and Fréchet (2011)</i>	0.961	>> 0.342	≈	0.307	>>	0.081	0.849	>>	0.324	≈	0.364	>>	0.054
<i>Dal Bó and Fréchet (2015)</i>	0.95	>> 0.265	≈	0.301	>>	0.081	0.843	>>	0.328	≈	0.369	>>	0.052
<i>Dreber et al. (2008)</i>	0.901	>> 0.154	≈	0.217	>>	0.062	1	>>	0.359	≈	0.203	>>	0.031
<i>Duffy and Ochs (2009)</i>	0.932	>> 0.218	≈	0.301	≈	0.208	0.748	>>	0.361	≈	0.35	>>	0.102
<i>Fréchet and Yuksel (2017)</i>	0.942	>> 0.132	≈	0.245	>>	0	1	≈	0.182	≈	0.364	≈	0.111
<i>Fudenberg et al. (2012)</i>	0.985	>> 0.429	≈	0.408	>>	0	0.921	>>	0.377	≈	0.443	>>	0.068
<i>Kagel and Schley (2013)</i>	0.947	>> 0.236	≈	0.288	>>	0.133	0.763	>>	0.298	≈	0.305	>>	0.042
<i>Sherstyuk et al. (2013)</i>	0.953	>> 0.312	≈	0.395	>>	0.172	0.875	>>	0.343	≈	0.349	>>	0.107
Pooled	0.949	>> 0.278	≈	0.3	>>	0.091	0.825	>>	0.333	≈	0.346	>>	0.059
Second halves per session													
<i>Aoyagi and Frechette (2009)</i>	0.965	>> 0.438	≈	0.625	≈	0.333	0.846	>>	0.371	<	0.443	≈	0.378
<i>Blonski et al. (2011)</i>	0.922	>> 0.157	≈	0.232	>>	0.027	0.941	>>	0.425	≈	0.23	>>	0.019
<i>Bruttel and Kamecke (2012)</i>	0.946	>> 0.156	≈	0.233	≈	0.173	0.958	>>	0.327	≈	0.4	>>	0.019
<i>Dal Bó (2005)</i>	0.918	>> 0.178	<	0.4	>	0.131	0.937	>>	0.32	≈	0.373	>>	0.052
<i>Dal Bó and Fréchet (2011)</i>	0.981	>> 0.373	≈	0.323	>>	0.077	0.95	>>	0.38	≈	0.416	>>	0.025
<i>Dal Bó and Fréchet (2015)</i>	0.98	>> 0.264	<	0.366	>>	0.058	0.904	>>	0.369	≈	0.44	>>	0.031
<i>Dreber et al. (2008)</i>	0.913	>> 0.029	≪	0.314	>>	0.007	0.955	>>	0.417	≈	0.611	>>	0.009
<i>Duffy and Ochs (2009)</i>	0.981	>> 0.362	≈	0.433	≈	0.226	0.889	>>	0.369	≈	0.368	>>	0.077
<i>Fréchet and Yuksel (2017)</i>	0.976	>> 0.173	≈	0.308	≈	0.222	0.75	>	0.294	≈	0.49	>>	0.06
<i>Fudenberg et al. (2012)</i>	0.976	>> 0.473	≈	0.509	≈	0.2	0.854	>>	0.5	≈	0.328	>>	0.077
<i>Kagel and Schley (2013)</i>	0.969	>> 0.218	≈	0.293	>	0.098	0.868	>>	0.332	≈	0.394	>>	0.02
<i>Sherstyuk et al. (2013)</i>	0.974	>> 0.465	≈	0.486	>>	0.107	0.952	>>	0.505	≈	0.369	>>	0.072
Pooled	0.975	>> 0.282	≪	0.351	>>	0.07	0.908	>>	0.378	≈	0.404	>>	0.033

Note: Relation signs, bootstrap procedure, and derived p -values are exactly as above, see Table 1, with the obvious adaptation that the Holm-Bonferroni correction now applies to all eight tests per data set.

Table 24: 1- and 2-memory SG behavior strategies versus best mixtures (by treatment) of 1- and 2-memory pure strategies (No switching)
(ICL-BIC of the models, less is better and relation signs point toward better models)

	SG+ SG M2“General”		SG M2“General”		Semi-Grim		Best Pure		Pure M1+G2,T2		Pure M1
Specification											
# Models evaluated	1		1		1		5		1		1
# Pars estimated (by treatment)	7		3		3		32		5		3
# Parameters accounted for	7		3		3		3-8		5		3
First halves per session											
<i>Aoyagi and Frechette (2009)</i>	743.26	≈	742.13	≫	694.72	<	763.74	≈	791.34	≈	800.02
<i>Blonski et al. (2011)</i>	713.93	≫	585.39	≫	549.45	≪	681.8	<	741.88	≫	707.94
<i>Bruttel and Kamecke (2012)</i>	569.53	≈	570.41	≈	567.86	≈	588.53	≈	588.53	≈	603.31
<i>Dal Bó (2005)</i>	418.65	≫	364.64	≈	358.51	<	388.99	≈	393.09	≈	394.31
<i>Dal Bó and Fréchet (2011)</i>	3576.79	≈	3594.64	≈	3533.99	≪	3837.37	≈	3836.86	<	3934.11
<i>Dal Bó and Fréchet (2015)</i>	5259.51	≫	5006.42	≈	4991.74	≪	5531.56	≈	5552.56	≈	5595.28
<i>Dreber et al. (2008)</i>	455.42	≈	451.44	≈	437.17	≈	462.66	≈	470.47	≈	462.66
<i>Duffy and Ochs (2009)</i>	1109.12	≈	1089.19	≈	1090.22	≈	1100.13	≈	1102.62	≈	1132.22
<i>Fréchet and Yuksel (2017)</i>	169.83	≈	164.74	≈	161.45	≪	181.95	≈	188.02	≈	181.95
<i>Fudenberg et al. (2012)</i>	313.22	≈	298.69	≈	291.43	≪	357.33	≈	366.77	≈	375.41
<i>Kagel and Schley (2013)</i>	1739.6	≈	1787.59	≈	1782.82	≈	1805.94	≈	1805.94	≈	1818.45
<i>Sherstyuk et al. (2013)</i>	918.44	≈	924.56	≈	912.8	≈	941.91	≈	941.91	≈	960.82
Pooled	16242.61	≫	15762.21	≫	15481.59	≪	16803.32	<	16962.38	≈	17075.91
Second halves per session											
<i>Aoyagi and Frechette (2009)</i>	437.99	≈	433.04	≫	389.24	≪	476.75	≈	484.39	≈	487.15
<i>Blonski et al. (2011)</i>	988.68	≫	888.46	≈	867.87	≪	990.66	≈	1084.43	≫	1067.21
<i>Bruttel and Kamecke (2012)</i>	333.93	≈	342.01	≈	347.4	≈	316.34	≈	316.34	≈	343.43
<i>Dal Bó (2005)</i>	475.66	≫	422.93	≈	424.44	<	461.93	≈	470.29	≈	464.4
<i>Dal Bó and Fréchet (2011)</i>	2758.96	<	2842.94	≈	2817.31	≈	2881.22	≈	2886.16	≪	3251.04
<i>Dal Bó and Fréchet (2015)</i>	5109.07	>	5027.27	≈	5043.81	≪	5564.25	≈	5586.2	≪	5853.42
<i>Dreber et al. (2008)</i>	271.57	≈	271.55	≈	264.94	≈	287.56	<	293.87	≈	287.56
<i>Duffy and Ochs (2009)</i>	1462.19	≈	1446.73	>	1403.03	≪	1617.75	≈	1617.75	≈	1661.55
<i>Fréchet and Yuksel (2017)</i>	348.73	≫	315.35	≈	313.5	<	356.1	≈	360.69	≈	356.1
<i>Fudenberg et al. (2012)</i>	386.55	≈	389.34	≈	380.75	<	445.16	≈	447.17	≈	476.66
<i>Kagel and Schley (2013)</i>	1161.5	≈	1208.7	≈	1211.37	≈	1169.29	≈	1169.29	<	1274.75
<i>Sherstyuk et al. (2013)</i>	551.82	<	596.62	≈	586.72	≈	583.79	≈	583.79	≪	691.06
Pooled	14541.97	>	14367.31	≫	14159.8	≪	15319.22	≈	15482.74	≪	16323.75

Note: Relation signs, bootstrap procedure, and derived p -values are exactly as above, see Table 1. Pure M1 refers to TFT, Grim, and AD. G2 denotes Grim2. For definitions of the strategies see Table 5.

Table 26: 1- and 2-memory SG behavior strategies versus best mixtures (by treatment) of 1- and 2-memory pure strategies (Random switching)
(ICL-BIC of the models, less is better and relation signs point toward better models)

	SG+SG M2“General”		SG M2“General”		Semi-Grim		Best Pure		Pure 1+G2,T2		Pure 1
Specification											
# Models evaluated	1		1		1		5		1		1
# Pars estimated (by treatment)	7		3		3		32		5		3
# Parameters accounted for	7		3		3		3-8		5		3
First halves per session											
<i>Aoyagi and Frechette (2009)</i>	742.43	≈	742.13	≫	694.72	<	763.63	≈	763.63	≈	793.71
<i>Blonski et al. (2011)</i>	703.59	≫	585.39	≫	549.45	≪	652.22	≪	691.76	>	659.7
<i>Bruttel and Kamecke (2012)</i>	584.62	≈	570.41	≈	567.86	≈	577.25	≈	577.25	≈	594.09
<i>Dal Bó (2005)</i>	376.06	>	364.64	≈	358.51	≈	354.58	≈	356.82	≈	365.42
<i>Dal Bó and Fréchette (2011)</i>	3731.81	≫	3594.64	≈	3533.99	≪	3820.5	≈	3827.21	<	3954
<i>Dal Bó and Fréchette (2015)</i>	5276.54	≫	5006.42	≈	4991.74	≪	5270.22	≈	5273.98	≪	5396.28
<i>Dreber et al. (2008)</i>	466.07	>	451.44	≈	437.17	≈	452.15	≈	452.15	≈	467.69
<i>Duffy and Ochs (2009)</i>	1109.16	≈	1089.19	≈	1090.22	≈	1100.74	≈	1100.74	≈	1119.78
<i>Fréchette and Yuksel (2017)</i>	201.85	≫	164.74	≈	161.45	≈	161.75	≈	164.76	≈	161.75
<i>Fudenberg et al. (2012)</i>	300.67	≈	298.69	≈	291.43	<	345.19	≈	345.19	<	371.76
<i>Kagel and Schley (2013)</i>	1805.7	≈	1787.59	≈	1782.82	≈	1788.69	≈	1788.69	<	1880.45
<i>Sherstyuk et al. (2013)</i>	949.23	≈	924.56	≈	912.8	<	958.93	≈	958.93	≈	977.38
Pooled	16503.06	≫	15762.21	≫	15481.59	≪	16405.97	≈	16483.49	≪	16851.43
Second halves per session											
<i>Aoyagi and Frechette (2009)</i>	433.73	≈	433.04	≫	389.24	<	458.42	≈	472.12	≈	475.15
<i>Blonski et al. (2011)</i>	982.16	≫	888.46	≈	867.87	<	1007.94	≈	1030.01	>	1013.07
<i>Bruttel and Kamecke (2012)</i>	340.43	≈	342.01	≈	347.4	≈	340.21	≈	340.21	≈	365
<i>Dal Bó (2005)</i>	431.83	≈	422.93	≈	424.44	≈	442.7	≈	444.23	≈	450.36
<i>Dal Bó and Fréchette (2011)</i>	2882.8	≈	2842.94	≈	2817.31	≪	3124.7	≈	3124.9	≪	3340.71
<i>Dal Bó and Fréchette (2015)</i>	5132.38	≫	5027.27	≈	5043.81	≪	5527	≈	5539.56	≪	5794.56
<i>Dreber et al. (2008)</i>	278.6	≈	271.55	≈	264.94	≈	285.14	≈	285.14	≈	285.64
<i>Duffy and Ochs (2009)</i>	1449.15	≈	1446.73	>	1403.03	≪	1587.69	≈	1587.69	<	1677.02
<i>Fréchette and Yuksel (2017)</i>	349.28	≫	315.35	≈	313.5	≈	334.45	≈	335.14	≈	334.45
<i>Fudenberg et al. (2012)</i>	394.42	≈	389.34	≈	380.75	≪	440.47	≈	440.47	<	488.69
<i>Kagel and Schley (2013)</i>	1204.73	≈	1208.7	≈	1211.37	≈	1240.91	≈	1240.91	<	1334.77
<i>Sherstyuk et al. (2013)</i>	587.54	≈	596.62	≈	586.72	<	658.71	≈	658.71	<	740.5
Pooled	14722.38	≫	14367.31	≫	14159.8	≪	15620.2	≈	15681.45	≪	16409.35

Note: Relation signs, bootstrap procedure, and derived p -values are exactly as above, see Table 1. Pure M1 refers to TFT, Grim, and AD. G2 denotes Grim2. For definitions of the strategies see Table 5.

Table 28: 1-memory or 2-memory Semi-Grim strategies, complexity of memory, mixtures of 1-memory and 2-memory SG (no switching)
(ICL-BIC of the models, less is better and relation signs point toward better models)

	SG M2“General”		SG M2“Semi-Grim”		SG M2“Grim”		Semi-Grim		SG M1 + M2“Grim”		SG M1 + M2“Gene
Specification											
# Models evaluated	1		1		1		1		1		1
# Pars estimated (by treatment)	5		4		3		3		5		7
# Parameters accounted for	5		4		3		3		5		7
First halves per session											
<i>Aoyagi and Frechette (2009)</i>	742.13	≈	740.41	≈	738.78	≫	694.72	≈	740.47	≈	743.27
<i>Blonski et al. (2011)</i>	585.39	≈	570.96	≫	551.67	≈	549.45	≪	681.49	≪	714.05
<i>Bruttel and Kamecke (2012)</i>	570.41	≈	568.94	≈	567.81	≈	567.86	≈	569.25	≈	569.55
<i>Dal Bó (2005)</i>	364.64	≈	360.06	≈	359.02	≈	358.51	≪	395.42	≪	418.7
<i>Dal Bó and Fréchette (2011)</i>	3594.64	≈	3588.23	≈	3577.91	≈	3533.99	≈	3557.12	≈	3576.96
<i>Dal Bó and Fréchette (2015)</i>	5006.42	≈	5011.68	≈	5002.07	≈	4991.74	≪	5225.84	≈	5259.73
<i>Dreber et al. (2008)</i>	451.44	≈	447.65	≈	444.6	≈	437.17	≈	454.47	≈	455.41
<i>Duffy and Ochs (2009)</i>	1089.19	≈	1088.56	≈	1087.11	≈	1090.22	≈	1109.13	≈	1109.16
<i>Fréchette and Yuksel (2017)</i>	164.74	≈	162.94	≈	161.74	≈	161.45	≈	165.9	≈	169.85
<i>Fudenberg et al. (2012)</i>	298.69	≈	300.14	≈	298.27	≈	291.43	<	310.29	≈	313.25
<i>Kagel and Schley (2013)</i>	1787.59	≈	1785.44	≈	1783.73	≈	1782.82	>	1736.56	≈	1739.6
<i>Sherstyuk et al. (2013)</i>	924.56	≈	925.17	≈	923.27	≈	912.8	≈	916.4	≈	918.43
Pooled	15762.21	≈	15696.08	>	15605.39	>	15481.59	≪	16044.72	≪	16243.29
Second halves per session											
<i>Aoyagi and Frechette (2009)</i>	433.04	≈	431.24	≈	430.13	>	389.24	<	436.18	≈	438.02
<i>Blonski et al. (2011)</i>	888.46	≈	889.64	>	879.23	>	867.87	<	962.97	≈	990.65
<i>Bruttel and Kamecke (2012)</i>	342.01	≈	340.72	≈	342.71	≈	347.4	≈	335.98	≈	333.93
<i>Dal Bó (2005)</i>	422.93	≈	419.05	≈	423.8	≈	424.44	≪	479.26	≈	475.76
<i>Dal Bó and Fréchette (2011)</i>	2842.94	≈	2844.4	≈	2835.16	≈	2817.31	≈	2745.37	≈	2759
<i>Dal Bó and Fréchette (2015)</i>	5027.27	≈	5059.55	≈	5058.67	≈	5043.81	≈	5128.18	≈	5109.14
<i>Dreber et al. (2008)</i>	271.55	≈	269.42	≈	266.02	≈	264.94	≈	267.81	≈	271.58
<i>Duffy and Ochs (2009)</i>	1446.73	≈	1444.42	≈	1442.23	≈	1403.03	≈	1461.69	≈	1462.25
<i>Fréchette and Yuksel (2017)</i>	315.35	≈	316.49	≈	314.7	≈	313.5	≪	333.42	<	348.67
<i>Fudenberg et al. (2012)</i>	389.34	≈	387.59	≈	385.81	≈	380.75	≈	384.43	≈	386.55
<i>Kagel and Schley (2013)</i>	1208.7	≈	1207.68	≈	1206.72	≈	1211.37	>	1156.83	≈	1161.49
<i>Sherstyuk et al. (2013)</i>	596.62	≈	596.25	≈	595.17	≈	586.72	≈	554.97	≈	551.82
Pooled	14367.31	≈	14352.35	≈	14289.77	>	14159.8	<	14429.45	≈	14544.19

Note: Relation signs, bootstrap procedure, and derived p -values are exactly as above.

Table 29: Table 28 by treatments – 1-memory or 2-memory Semi-Grim strategies, complexity of memory, mixtures of 1-memory and 2-memory SG (no switching)

(a) First halves per session

(b) Second halves per session

	SG M2"General"	SG M2"Semi-Grim"	SG M2"Grim"	Semi-Grim	SG M1 + M2"Grim"	SG M1 + M2"General"
Specification						
# Models evaluated	1	1	1	1	1	1
# Pars estimated (by treatment)	5	4	3	3	5	7
# Parameters accounted for	5	4	3	3	5	7
AF09–34	742.13	740.41	738.78	694.72	740.47	743.27
BOS11–9	29.54	28.04	28.15	28.34	43.69	44.2
BOS11–14	45.65	44.8	43.39	43.74	56.61	51.98
BOS11–15	17.23	15.73	14.23	14.27	30.2	33.19
BOS11–16	102.32	101.02	100.59	102.31	113.78	114.86
BOS11–17	54.59	53.09	51.86	48.33	68.7	71.21
BOS11–26	136.72	135.26	133.94	130.46	157.04	158.29
BOS11–27	46.89	46.76	46.18	49.15	53.4	55.81
BOS11–30	19.4	17.9	16.4	15.97	33.26	36.25
BOS11–31	82.99	88.31	86.87	86.84	74.74	78.15
BK12–28	570.41	568.94	567.81	567.86	569.25	569.55
D05–18	140.24	138.37	136.52	136.21	165.24	161.88
D05–19	220.85	218.85	220.37	220.17	226.64	251.86
DF11–6	337.83	339.19	337.3	333.68	339.43	346.41
DF11–7	525.23	528.44	527.32	512.99	526.64	530.12
DF11–8	557.31	557.06	555.25	551.81	524.93	525.4
DF11–22	685.45	683.66	686.66	684.04	691.45	693.27
DF11–23	664.99	663.28	661.46	656.29	657.42	653.29
DF11–24	796.91	795.06	793.78	779.03	790.32	790.78
DF15–4	160.51	158.66	157.78	149.71	186.25	189.64
DF15–5	716.06	714.63	712.46	716.51	780.89	785.3
DF15–20	916.66	917.24	916.3	914.2	963.41	955.65
DF15–21	1201.29	1198.83	1196.6	1196.29	1211.55	1228.04
DF15–33	1695.1	1711.9	1716.1	1714.32	1763.72	1766.98
DF15–35	287.64	287.1	285.35	283.21	290.86	293.3
DRFN08–10	191.06	189.47	187.8	180.96	198.75	195.66
DRFN08–11	256.88	255.38	254.69	254.11	252.22	254.85
DO09–32	1089.19	1088.56	1087.11	1090.22	1109.13	1109.16
FY17–25	164.74	162.94	161.74	161.45	165.9	169.85
FRD12–29	298.69	300.14	298.27	291.43	310.29	313.25
KS13–12	1787.59	1785.44	1783.73	1782.82	1736.56	1739.6
STS13–13	924.56	925.17	923.27	912.8	916.4	918.43
<i>Aoyagi and Fréchette (2009)</i>	742.13	740.41	738.78	694.72	740.47	743.27
<i>Blonski et al. (2011)</i>	585.39	570.96	551.67	549.45	681.49	714.05
<i>Brützel and Kamecke (2012)</i>	570.41	568.94	567.81	567.86	569.25	569.55
<i>Dal Bó (2005)</i>	364.64	360.06	359.02	358.51	395.42	418.7
<i>Dal Bó and Fréchette (2011)</i>	3594.64	3588.23	3577.91	3533.99	3557.12	3576.96
<i>Dal Bó and Fréchette (2015)</i>	5006.42	5011.68	5002.07	4991.74	5225.84	5259.73
<i>Dreber et al. (2008)</i>	451.44	447.65	444.6	437.17	454.47	455.41
<i>Duffy and Ochs (2009)</i>	1089.19	1088.56	1087.11	1090.22	1109.13	1109.16
<i>Fréchette and Yuksel (2017)</i>	164.74	162.94	161.74	161.45	165.9	169.85
<i>Fudenberg et al. (2012)</i>	298.69	300.14	298.27	291.43	310.29	313.25
<i>Kagel and Schley (2013)</i>	1787.59	1785.44	1783.73	1782.82	1736.56	1739.6
<i>Sherstyuk et al. (2013)</i>	924.56	925.17	923.27	912.8	916.4	918.43
Pooled	15762.21	15696.08	15605.39	15481.59	16044.72	16243.29

	SG M2"General"	SG M2"Semi-Grim"	SG M2"Grim"	Semi-Grim	SG M1 + M2"Grim"	SG M1 + M2"General"
Specification						
# Models evaluated	1	1	1	1	1	1
# Pars estimated (by treatment)	5	4	3	3	5	7
# Parameters accounted for	5	4	3	3	5	7
AF09–34	1181.85	1180.06	1178.95	1085.04	1179.77	1181.23
BOS11–9	51.41	50.6	50.3	50.83	64.15	62.29
BOS11–14	63.05	62.16	61.57	61.74	78.27	76.73
BOS11–15	27.05	25.55	24.05	21.88	40.08	43.07
BOS11–16	180.95	179.45	178.66	178.84	185.48	185.96
BOS11–17	192.93	191.5	190	183.03	205.05	195.93
BOS11–26	319.54	317.83	318.19	317.33	331.68	333.84
BOS11–27	214.41	229.34	235.12	231.72	217.95	212.66
BOS11–30	94.64	93.32	93.07	93.22	105.76	108.3
BOS11–31	185.95	198.29	200.1	200.87	166.11	166.7
BK12–28	915.39	914.36	915.95	916.81	895.53	891.9
D05–18	356.04	354.4	352.82	351.89	378.37	368.38
D05–19	409.37	407.41	417.01	418.28	420.37	431.08
DF11–6	590.39	596.23	595.49	599.29	579.49	582.07
DF11–7	1197.17	1205.8	1204.41	1199.27	1153.19	1151.03
DF11–8	996.09	999.9	998.05	985.67	915.18	924.46
DF11–22	1327.01	1326.07	1332.1	1331.85	1293.25	1294.96
DF11–23	1135.01	1133.45	1131.71	1113.94	1099.02	1099.93
DF11–24	1159.28	1157.61	1156.21	1122.26	1161.89	1163.33
DF15–4	271.14	269.69	269.21	260.02	296.55	297.48
DF15–5	1323.6	1324.74	1332.34	1331.49	1332.32	1336.61
DF15–20	1706.95	1714.68	1717.74	1716.02	1691.17	1696.55
DF15–21	1223.11	1223.36	1229.41	1223.88	2182.89	2186.44
DF15–33	3612.75	3661.3	3679	3678.61	3682.31	3618.49
DF15–35	858.27	862.75	863.47	856.78	845.23	845.69
DRFN08–10	289.36	287.86	286.23	270.35	292.64	295.95
DRFN08–11	426.29	425.3	424.52	425.89	411.22	413.33
DO09–32	2571.83	2570.1	2568.54	2526.02	2557.19	2559.18
FY17–25	473.7	474.86	473.42	473.84	490.2	497.64
FRD12–29	680.53	681.05	679.32	668.18	656.01	658.77
KS13–12	3005.9	3004.83	3004.85	3009.28	2876.28	2880.12
STS13–13	1528.33	1530.53	1529.48	1507.33	1469.2	1472.88
<i>Aoyagi and Fréchette (2009)</i>	1181.85	1180.06	1178.95	1085.04	1179.77	1181.23
<i>Blonski et al. (2011)</i>	1380.01	1388.1	1381.12	1369.51	1444.6	1454.59
<i>Brützel and Kamecke (2012)</i>	915.39	914.36	915.95	916.81	895.53	891.9
<i>Dal Bó (2005)</i>	768.96	764.65	771.96	772.3	802.29	804.41
<i>Dal Bó and Fréchette (2011)</i>	6431.87	6440.62	6434.12	6368.43	6228.94	6253.48
<i>Dal Bó and Fréchette (2015)</i>	10034.99	10087.85	10099.67	10074.29	10059.64	10022.08
<i>Dreber et al. (2008)</i>	719.15	715.96	712.86	698.35	707.36	714.18
<i>Duffy and Ochs (2009)</i>	2571.83	2570.1	2568.54	2526.02	2557.19	2559.18
<i>Fréchette and Yuksel (2017)</i>	473.7	474.86	473.42	473.84	490.2	497.64
<i>Fudenberg et al. (2012)</i>	680.53	681.05	679.32	668.18	656.01	658.77
<i>Kagel and Schley (2013)</i>	3005.9	3004.83	3004.85	3009.28	2876.28	2880.12
<i>Sherstyuk et al. (2013)</i>	1528.33	1530.53	1529.48	1507.33	1469.2	1472.88
Pooled	29874.89	29898.86	29859.67	29578.81	29549.38	29645.79

Note: Notation of treatments and meaning of relation signs are all as defined above, see Table 7.

Table 30: Mixtures of 1- and 2-memory pure and generalized strategies (no switching)
(ICL-BIC of the models, less is better and relation signs point toward better models)

	Gen M2		Gen M1		Best Pure M2		Pure M1		+ G2, TFT2, T2		+ 2TFT	
Specification												
# Models evaluated	1		1		5		1		1		1	
# Pars estimated (by treatment)	9		6		32		3		6		7	
# Parameters accounted for	9		6		3–8		3		6		7	
First halves per session												
<i>Aoyagi and Frechette (2009)</i>	649.81	≈	645.32	≪	791.39	≈	800.02	≈	797.53	≈	799.33	
<i>Blonski et al. (2011)</i>	767.03	≈	748.97	≈	705.79	≈	708.13	≪	764.55	≪	795.15	
<i>Bruttel and Kamecke (2012)</i>	581.45	≈	617.33	≈	588.54	≈	603.32	≈	590.34	≈	592.12	
<i>Dal Bó (2005)</i>	402.78	≈	407.86	≈	389.07	≈	394.44	≈	398.94	≪	414.37	
<i>Dal Bó and Fréchette (2011)</i>	3577.56	≈	3624.47	≪	3835.49	<	3934.17	≈	3848.71	≈	3866.15	
<i>Dal Bó and Fréchette (2015)</i>	5349.43	≈	5355.13	≪	5538.04	≈	5595.6	≈	5573.66	≈	5609.72	
<i>Dreber et al. (2008)</i>	487.89	≈	479.28	>	462.71	≈	462.71	≈	473.97	≈	479.06	
<i>Duffy and Ochs (2009)</i>	1059.59	≈	1056.99	≈	1102.65	≈	1132.23	≈	1106.99	≈	1116.34	
<i>Fréchette and Yuksel (2017)</i>	188.15	≈	188.5	≈	181.98	≈	181.98	≪	195.97	<	202.98	
<i>Fudenberg et al. (2012)</i>	322.89	≈	324.36	<	366.75	≈	375.41	≈	368.71	≈	370.66	
<i>Kagel and Schley (2013)</i>	1684.63	<	1779.94	≈	1805.95	≈	1818.46	≈	1808.32	≈	1812.67	
<i>Sherstyuk et al. (2013)</i>	896.24	≈	903.93	≈	941.92	≈	960.83	≈	943.92	≈	946.88	
Pooled	16295.72	≈	16350.93	≪	16862.01	<	17076.73	≈	17090.45	<	17260.76	
Second halves per session												
<i>Aoyagi and Frechette (2009)</i>	365.6	≈	363.58	≪	484.4	≈	487.16	≈	490.13	≈	488.8	
<i>Blonski et al. (2011)</i>	1022.73	≈	1094.17	≈	1056	≈	1067.5	≈	1108.59	≈	1147.6	
<i>Bruttel and Kamecke (2012)</i>	340.71	≈	348.23	≈	316.38	≈	343.47	≈	318.17	≈	319.93	
<i>Dal Bó (2005)</i>	471.65	≈	475.11	≈	463.52	≈	464.42	≈	476.02	≈	483.87	
<i>Dal Bó and Fréchette (2011)</i>	2706.97	≪	2988.87	≈	2880.23	≪	3251.2	≫	2903.67	≈	2914.93	
<i>Dal Bó and Fréchette (2015)</i>	5227.69	<	5448.95	≈	5577.08	≪	5853.69	≫	5603.14	≈	5635.29	
<i>Dreber et al. (2008)</i>	288.95	≈	295.06	≈	287.58	≈	287.58	<	298.37	≈	301.17	
<i>Duffy and Ochs (2009)</i>	1370.98	≈	1409.56	≪	1617.77	≈	1661.56	≈	1620.05	≈	1622.69	
<i>Fréchette and Yuksel (2017)</i>	311.99	≈	309.63	≪	356.11	≈	356.11	≈	362.81	≈	364.77	
<i>Fudenberg et al. (2012)</i>	364.47	≈	373.44	<	447.19	≈	476.66	>	449.08	≈	449.92	
<i>Kagel and Schley (2013)</i>	1128.24	≈	1211.25	≈	1169.31	<	1274.78	>	1171.67	≈	1172	
<i>Sherstyuk et al. (2013)</i>	551.88	<	616.19	≈	583.8	≪	691.06	≫	585.8	≈	587.81	
Pooled	14480.14	≪	15152.88	≈	15396.72	≪	16324.6	≫	15606.38	≈	15744.12	

Note: Relation signs, bootstrap procedure, and derived p -values are exactly as above, see Table 1. Pure M1 refers to TFT, Grim, and AD. G2 denotes Grim2. For definitions of pure strategies see Table 5. Gen M1 refers to generalized versions of TFT, Grim, and AD with memory-1. “+ G2, TFT2, T2” adds those strategies to the set of “Pure M1”. “+2TFT” adds this strategy on top of the former.

Table 32: Comparison of 1- and 2-memory Semi-Grim with two and three parameters, pure and generalized strategies (no switching, Grim scheme)
 (ICL-BIC of the models, less is better and relation signs point toward better models)

	SGs M2“General”		SGs M2 “Grim”		Semi-Grim		Gen M2“Grim”		Gen M1		Best Pure M2
Specification											
# Models evaluated	1		1		1		1		1		5
# Pars estimated (by treatment)	5		3		3		9		6		32
# Parameters accounted for	5		3		3		9		6		3–8
First halves per session											
<i>Aoyagi and Frechette (2009)</i>	742.13	≈	738.78	≫	694.72	>	649.81	≈	645.32	≪	791.39
<i>Blonski et al. (2011)</i>	585.39	≫	551.67	≈	549.45	≪	767.03	≈	748.97	≈	705.79
<i>Bruttel and Kamecke (2012)</i>	570.41	≈	567.81	≈	567.86	≈	581.45	≈	617.33	≈	588.54
<i>Dal Bó (2005)</i>	364.64	≈	359.02	≈	358.51	<	402.78	≈	407.86	≈	389.07
<i>Dal Bó and Fréchette (2011)</i>	3594.64	≈	3577.91	≈	3533.99	≈	3577.56	≈	3624.47	≪	3835.49
<i>Dal Bó and Fréchette (2015)</i>	5006.42	≈	5002.07	≈	4991.74	≪	5349.43	≈	5355.13	≪	5538.04
<i>Dreber et al. (2008)</i>	451.44	≈	444.6	≈	437.17	<	487.89	≈	479.28	>	462.71
<i>Duffy and Ochs (2009)</i>	1089.19	≈	1087.11	≈	1090.22	≈	1059.59	≈	1056.99	≈	1102.65
<i>Fréchette and Yuksel (2017)</i>	164.74	≈	161.74	≈	161.45	≪	188.15	≈	188.5	≈	181.98
<i>Fudenberg et al. (2012)</i>	298.69	≈	298.27	≈	291.43	<	322.89	≈	324.36	<	366.75
<i>Kagel and Schley (2013)</i>	1787.59	≈	1783.73	≈	1782.82	>	1684.63	<	1779.94	≈	1805.95
<i>Sherstyuk et al. (2013)</i>	924.56	≈	923.27	≈	912.8	≈	896.24	≈	903.93	≈	941.92
Pooled	15762.21	≫	15605.39	>	15481.59	≪	16295.72	≈	16350.93	≪	16862.01
Second halves per session											
<i>Aoyagi and Frechette (2009)</i>	433.04	≈	430.13	>	389.24	≈	365.6	≈	363.58	≪	484.4
<i>Blonski et al. (2011)</i>	888.46	≈	879.23	>	867.87	≪	1022.73	≈	1094.17	≈	1056
<i>Bruttel and Kamecke (2012)</i>	342.01	≈	342.71	≈	347.4	≈	340.71	≈	348.23	≈	316.38
<i>Dal Bó (2005)</i>	422.93	≈	423.8	≈	424.44	≪	471.65	≈	475.11	≈	463.52
<i>Dal Bó and Fréchette (2011)</i>	2842.94	≈	2835.16	≈	2817.31	≈	2706.97	≪	2988.87	≈	2880.23
<i>Dal Bó and Fréchette (2015)</i>	5027.27	≈	5058.67	≈	5043.81	<	5227.69	<	5448.95	≈	5577.08
<i>Dreber et al. (2008)</i>	271.55	≈	266.02	≈	264.94	≈	288.95	≈	295.06	≈	287.58
<i>Duffy and Ochs (2009)</i>	1446.73	≈	1442.23	≈	1403.03	≈	1370.98	≈	1409.56	≪	1617.77
<i>Fréchette and Yuksel (2017)</i>	315.35	≈	314.7	≈	313.5	≈	311.99	≈	309.63	≪	356.11
<i>Fudenberg et al. (2012)</i>	389.34	≈	385.81	≈	380.75	≈	364.47	≈	373.44	<	447.19
<i>Kagel and Schley (2013)</i>	1208.7	≈	1206.72	≈	1211.37	>	1128.24	≈	1211.25	≈	1169.31
<i>Sherstyuk et al. (2013)</i>	596.62	≈	595.17	≈	586.72	≈	551.88	<	616.19	≈	583.8
Pooled	14367.31	≈	14289.77	>	14159.8	≈	14480.14	≪	15152.88	≈	15396.72

Note: Relation signs, bootstrap procedure, and derived p -values are exactly as above, see Table 1. Pure M1 refers to TFT, Grim, and AD. For definitions of pure strategies see Table 5. Gen M1 refers to generalized versions of TFT, Grim, and AD with memory-1. SGs refers to a two parameter version of SG ($1 - \theta_1, \theta_2, \theta_2, \theta_1$). “Gen M2” refers to memory-2 versions of the generalized strategies that allow parameters to depend on the prevalence of joint cooperation in $t - 2$ (Grim Scheme).

Table 33: Comparison of 1- and 2-memory Semi-Grim, pure and generalized strategies (no switching, TFT scheme) (ICL-BIC of the models, less is better and relation signs point toward better models)

	SGs M2“General”		SGs M2 “TFT”		Semi-Grim		Gen M2“TFT”		Gen M1		Best Pure M2
Specification											
# Models evaluated	1		1		1		1		1		5
# Pars estimated (by treatment)	5		3		3		9		6		32
# Parameters accounted for	5		3		3		9		6		3–8
First halves per session											
<i>Aoyagi and Frechette (2009)</i>	742.13	≈	738.55	≫	694.72	>	649.37	≈	645.32	≪	791.38
<i>Blonski et al. (2011)</i>	585.39	≫	550.9	≈	549.45	≪	760.49	≈	748.97	≈	703.1
<i>Bruttel and Kamecke (2012)</i>	570.41	≈	567.91	≈	567.86	≈	580.92	≈	617.33	≈	588.55
<i>Dal Bó (2005)</i>	364.64	≈	359.74	≈	358.51	≪	405.83	≈	407.86	≈	389.08
<i>Dal Bó and Fréchette (2011)</i>	3594.64	≈	3573.88	≈	3533.99	≈	3524.12	<	3624.47	≪	3835.75
<i>Dal Bó and Fréchette (2015)</i>	5006.42	≈	5004.6	≈	4991.74	≪	5276.75	≈	5355.13	≪	5538.37
<i>Dreber et al. (2008)</i>	451.44	≈	445.38	≈	437.17	<	482.83	≈	479.28	>	462.71
<i>Duffy and Ochs (2009)</i>	1089.19	≈	1090.75	≈	1090.22	≈	1039.25	≈	1056.99	≈	1102.64
<i>Fréchette and Yuksel (2017)</i>	164.74	≈	161.91	≈	161.45	≪	182.81	≈	188.5	≈	181.98
<i>Fudenberg et al. (2012)</i>	298.69	≈	296.73	≈	291.43	<	319.76	≈	324.36	<	366.78
<i>Kagel and Schley (2013)</i>	1787.59	≈	1784.29	≈	1782.82	>	1651.6	≪	1779.94	≈	1805.95
<i>Sherstyuk et al. (2013)</i>	924.56	≈	922.74	≈	912.8	≈	890.54	≈	903.93	≈	941.91
Pooled	15762.21	≫	15606.79	>	15481.59	≪	16092.54	<	16350.93	≪	16858.86
Second halves per session											
<i>Aoyagi and Frechette (2009)</i>	433.04	≈	431.45	≫	389.24	≈	368.09	≈	363.58	≪	484.41
<i>Blonski et al. (2011)</i>	888.46	≈	874.57	≈	867.87	≪	1013.57	≈	1094.17	≈	1055.95
<i>Bruttel and Kamecke (2012)</i>	342.01	≈	345.65	≈	347.4	≈	324.09	≈	348.23	≈	316.38
<i>Dal Bó (2005)</i>	422.93	≈	424.63	≈	424.44	<	451.34	<	475.11	≈	463.54
<i>Dal Bó and Fréchette (2011)</i>	2842.94	≈	2833.79	≈	2817.31	≈	2679.96	≪	2988.87	≈	2885.43
<i>Dal Bó and Fréchette (2015)</i>	5027.27	≈	5048.25	≈	5043.81	≈	5144.02	≪	5448.95	≈	5577.55
<i>Dreber et al. (2008)</i>	271.55	≈	270.81	≈	264.94	≈	287.55	≈	295.06	≈	287.58
<i>Duffy and Ochs (2009)</i>	1446.73	≈	1442.34	≈	1403.03	≈	1349.18	<	1409.56	≪	1617.77
<i>Fréchette and Yuksel (2017)</i>	315.35	≈	312.94	≈	313.5	≈	311.04	≈	309.63	≪	356.11
<i>Fudenberg et al. (2012)</i>	389.34	≈	386.5	≈	380.75	≈	367.2	≈	373.44	<	447.19
<i>Kagel and Schley (2013)</i>	1208.7	≈	1208.4	≈	1211.37	>	1100.94	≈	1211.25	≈	1169.31
<i>Sherstyuk et al. (2013)</i>	596.62	≈	595.53	≈	586.72	>	541.54	<	616.19	≈	583.8
Pooled	14367.31	≈	14284.29	>	14159.8	≈	14266.78	≪	15152.88	≈	15402.36

Note: Relation signs, bootstrap procedure, and derived p -values are exactly as above, see Table 1. Pure M1 refers to TFT, Grim, and AD. For definitions of pure strategies see Table 5. “Gen M1” refers to generalized versions of TFT, Grim, and AD with memory-1. SGs refers to a two parameter version of SG ($1 - \theta_1, \theta_2, \theta_2, \theta_1$). “Gen M2” refers to memory-2 versions of the generalized strategies that allow parameters to depend on opponent’s behavior in $t - 2$ (TFT Scheme).

Table 34: Examining all mixtures of Semi-Grim with pure or generalized pure strategies as secondary components (robustness check for Table 4)

Component 1	First component is always Semi-Grim														
Component 2	Gen WSLs	Gen TFT	Gen Grim	Gen AD/AC	AD	Grim	TFT	WSLS							
Specification															
# Models evaluated	1	1	1	1	1	1	1	1	1	1	1	1			
# Pars estimated (by treatment)	5	5	5	5	5	4	4	4	4	4	4	4			
# Parameters accounted for	5	5	5	5	5	4	4	4	4	4	4	4			
First halves per session															
<i>Aoyagi and Fréchette (2009)</i>	684.3	≈	689.91	≈	688.17	≈	649.86	≈	698.36	≈	681.99	≈	698.37	≈	698.36
<i>Blonski et al. (2011)</i>	635.26	≈	648.6	≈	645.78	≫	614.35	≪	663.28	≈	658.75	>	614.21	>	594.32
<i>Bruttel and Kamecke (2012)</i>	572.36	≈	559.64	≈	571.61	≈	562.11	≈	567.84	≈	562.35	≈	571.1	≈	566.52
<i>Dal Bó (2005)</i>	393.92	≈	372.81	≈	384.74	≈	372.89	≈	378.87	≈	389.65	≈	386.98	>	368.21
<i>Dal Bó and Fréchette (2011)</i>	3538.98	≈	3532.69	≈	3475.81	≈	3416.01	≈	3488.81	≈	3434.87	<	3538.27	≈	3577.29
<i>Dal Bó and Fréchette (2015)</i>	5164.89	≫	5070.81	≈	5052.96	≫	4941.84	≈	5024.63	≈	5107.87	>	5047.74	≈	5033.06
<i>Dreber et al. (2008)</i>	454.45	≈	453.19	≈	440.94	≈	443.12	≈	442.47	≈	440.24	≈	443.9	≈	444.99
<i>Duffy and Ochs (2009)</i>	1084.07	≈	1062.36	≈	1074.19	>	1017.65	<	1084.74	≈	1074.22	≈	1099.57	≈	1084.81
<i>Fréchette and Yuksel (2017)</i>	180.17	≈	190.03	≈	188.6	≫	167.94	≈	165.94	<	178.18	≈	180.54	≫	165.36
<i>Fudenberg et al. (2012)</i>	296.17	≈	299.94	≈	299.98	≈	291.68	≈	292.47	≈	295.53	≈	297.99	≈	295.46
<i>Kagel and Schley (2013)</i>	1785.45	≈	1780	≈	1777.41	>	1690.55	≈	1689.17	≈	1685.56	<	1777.58	≈	1787.56
<i>Sherstyuk et al. (2013)</i>	910.79	≈	901.32	≈	911	>	866.46	≈	887.79	≈	889.32	≈	904.29	≈	916.83
Pooled	15919.66	>	15780.15	≈	15730.04	≫	15253.3	≪	15566.74	≈	15580.91	≈	15742.91	≈	15715.14
Second halves per session															
<i>Aoyagi and Fréchette (2009)</i>	395.39	≈	380.72	≈	403.19	>	357.21	≈	392.87	≈	398.09	≈	397.94	≈	393.02
<i>Blonski et al. (2011)</i>	921.84	≈	946.43	<	1157.9	≫	881.73	≈	914.98	≈	915.1	≈	940.75	>	919.01
<i>Bruttel and Kamecke (2012)</i>	323.04	≈	339	≈	337.37	≈	346.24	≈	349.32	≈	336.68	≈	337.22	≈	350.98
<i>Dal Bó (2005)</i>	471.52	≈	461.35	≈	444.45	≈	437.05	≈	437.71	<	457.32	≈	444.73	≈	433.69
<i>Dal Bó and Fréchette (2011)</i>	2729.84	≈	2742.23	≪	3652.93	≫	2620.29	≈	2640.88	≈	2615.25	<	2798.46	≈	2845.99
<i>Dal Bó and Fréchette (2015)</i>	5043.7	≈	5080.21	≈	5025.85	>	4925.77	≈	5032.59	≈	5077.82	≈	5053.74	≈	5042.75
<i>Dreber et al. (2008)</i>	279.59	≈	258.73	≈	276.69	≈	279.28	≈	275.84	≈	277.77	≈	262.93	≈	275.37
<i>Duffy and Ochs (2009)</i>	1407.13	≈	1378.8	≈	1365.2	>	1294.4	<	1399.65	≈	1379.73	≈	1414.84	≈	1447.23
<i>Fréchette and Yuksel (2017)</i>	313.59	≈	324.17	≈	323.61	>	283.31	<	319.78	≈	311.64	≈	319.44	≈	317.41
<i>Fudenberg et al. (2012)</i>	381.7	≈	387.98	≈	385.9	≈	364.37	≈	366.79	≈	379.95	≈	388.9	≈	384.62
<i>Kagel and Schley (2013)</i>	1156.16	≈	1202.31	≈	1202.28	≈	1154.97	≈	1126.89	≈	1153.79	≈	1200.87	≈	1216.1
<i>Sherstyuk et al. (2013)</i>	588.67	≈	569.49	≈	583.54	≈	565.48	≈	550.64	≈	539.35	≈	567.48	≈	594.92
Pooled	14231.02	≈	14290.27	≪	15377.76	≫	13728.94	<	13990.32	≈	14024.87	<	14309.69	≈	14403.47

Note: Relation signs, bootstrap procedure, and derived p -values are exactly as above, see Table 1. For definitions of pure strategies see Table 5. For definitions of generalized strategies see Section 3 main text.