# WORKING PAPER SERIES

## Dynamic Psychological  Games

*Pierpaolo Battigalli and Martin Dufwenberg*

**Working Paper n. 287**

April 2005

# Dynamic Psychological Games[*]

Pierpaolo Battigalli     Martin Dufwenberg

Bocconi University     University of Arizona

April 2005

## Abstract

Building on recent work on dynamic interactive epistemology, we extend the analysis of extensive-form psychological games (Geneakoplos, Pearce & Stacchetti, *Games and Economic Behavior*, 1989) to include conditional higher-order beliefs and enlarged domains of payoff functions. The approach allows modeling dynamic psychological effects (such as sequential reciprocity, psychological forward induction, and regret) that are ruled out when epistemic types are identified with hierarchies of initial beliefs. We define a notion of *psychological sequential equilibrium*, which generalizes the sequential equilibrium notion for traditional games, for which we prove existence under mild assumptions. Our framework also allows us to directly formulate assumptions about 'dynamic' rationality and interactive beliefs in order to explore strategic interaction without assuming that players beliefs are coordinated on an equilibrium. In particular, we provide an exploration of (extensive-form) rationalizability in psychological games.

KEYWORDS: psychological games, belief-dependent motivation, extensive-form solution concepts, dynamic interactive epistemology.

J.E.L. *classification numbers:* C72, C73.

# 1   Introduction

We develop a framework for analyzing strategic interaction when players have 'belief-dependent' motivations, thereby generalizing the theory of extensive form *psychological games* proposed by Geanakoplos, Pearce & Stacchetti (1989; henceforth GPS). The rest of this introduction motivates our approach in more detail.

Traditional game theory is not a rich enough toolbox to adequately describe many psychological or social aspects of motivation and behavior. The traditional approach assumes that payoffs depend only on which actions are chosen. By contrast, the payoffs of decision makers who are emotional or motivated by reciprocity or social respect may depend also on which beliefs (about choices, beliefs, or information) players harbor. The following examples illustrate:

1. When Ann takes a taxi ride she tips as much as she expects that the driver (Bob) expects to get. She suffers from guilt if she tips less.

2. Cleo suddenly pushes Dan over. Should Dan splash a bucket of water over Cleo in return? Maybe she actually tried to hug him? If so, Dan would rather forgive (maybe even hug) Cleo.

3. Eva is unemployed. Her neighbor, Fred, observes the effort with which she tries to get a job. Fred's taxes pay for Eva's unemployment benefits, so Eva's choice has externalities the size of which depends on her talent translating effort to probability of getting a job (low effort is costlier to Fred if Eva is talented and could have gotten a job had she tried harder). Eva's talent is known only to her, but Fred makes inferences observing her effort. This determines the social respect he bestows on Eva, and since she cares about respect this influences her effort.

Ann's tip, Dan's hug/soak choice, and Eva's effort each pins down a strategy profile. Yet the preferred choice depends on the player's belief.[1]

The point that belief-dependent motivation may be important for strategic decision making is made by GPS, who present several intriguing examples involving various emotions. They show the inadequacy of traditional methods to represent the involved preferences, and develop an extension (in the normal as well as in the extensive form) of traditional game theory to deal with the matter.[2] Only recently, however, has a larger set of economists

---

[1]Ann's preference depends on her belief of Bob's belief; Dan's on his assessment of Cleo's intentions; Eva's preferences over effort depend on Fred's inferences on her talent.

[2]Gilboa & Schmeidler (1988) also consider some games with belief-dependent payoffs.

come to acknowledge the relevance of belief-dependent motivation, mainly following the work by experimentalists.[3] In the lab, subjects often display 'non-selfish' behavior. This has inspired theoretical models of 'social preferences' which can rationalize the data.[4] These models differ in structure, and some do not require a deviation from traditional game theory (*e.g.*, inequity aversion models). However, a few models describe belief-dependent motivation, and some experiments support such models.[5]

While GPS' paper is highly inspiring for all this work, a careful scrutiny reveals that their approach is too restrictive to handle many plausible forms of belief-dependent motivation (this is acknowledged by GPS themselves; see pp. 70, 78). There are several reasons:

**R1 (updated beliefs):** GPS only allow *initial* beliefs to enter the domain of a player's utility, while many seemingly important forms of belief-dependent motivation require *updated* beliefs to matter.

**R2 (others' beliefs):** GPS only allow a player's *own* beliefs to enter the domain of his utility function, while there are conceptual and technical reasons to let *others'* beliefs matter.

**R3 (dependence on strategies):** GPS follow the traditional extensive games approach of letting strategies influence utilities only insofar as they influence terminal histories, but many forms of belief-dependent motivation become compelling in particular in conjunction with preferences that depend on strategies in ways not captured by terminal histories.

**R4 (non-equilibrium analysis):** GPS restrict attention to equilibrium analysis, but in many strategic situations there is little compelling reason to expect players to coordinate on an equilibrium and one may wish to explore alternative assumptions.

This list deserves more discussion and backup by examples, but we postpone this until the next section. Here we just note that items in the list

---

[3]See, however, the applied psychological-game theoretical work by Huang & Wu (1994), Dufwenberg (1995), Geanakoplos (1996), Ruffle (1999), Huck & Kübler (2000), Dufwenberg (2002), and Li (2005), as well as the models by Bernheim (1994) and Dufwenberg & Lundholm (2000) which can be given psychological-game interpretations (as we explain below) although such connections are not made in the original papers.

[4]See Fehr & Gächter (2000) for a discussion.

[5]For models, see Rabin (1993), Dufwenberg & Kirchsteiger (2004), Falk & Fischbacher (1998), and Charness & Dufwenberg (2004); for experiments, see Dufwenberg & Gneezy (2000), Bacharach, Guerra & Zizzo (2002), Guerra & Zizzo (2004), and Charness & Dufwenberg (2004).

have lead some researchers to deviate from GPS' framework, in developing specific examples or models where belief-dependent motivation play a role. However, apart from GPS, almost no papers are concerned with developing the overall framework of psychological game theory, by defining new classes of psychological games for which solution concepts are provided.[6] In this paper we attempt to fill this gap.

Our approach crucially draws on results and insights by Battigalli & Siniscalchi (1999) on how to represent hierarchies of conditional beliefs. This material is essential for a systematic treatment of **R1**, and figures in the background of **R2**-**R4** since these issues are relevant in contexts with updated beliefs. We define a large class of psychological games, which contains (in a particular sense) GPS' games as well as traditional games as special cases. We introduce a new notion of psychological sequential equilibrium, which generalizes Kreps & Wilson's (1982) sequential equilibrium notion, for which we prove an existence theorem. We develop a framework for analyzing interactive epistemology in psychological games, which we employ to develop a notion of rationalizability, thereby addressing **R4**.

Section 2 gives an overview of the conceptual issues that lie behind and motivate our work. Section 3 develops the general framework, and in particular defines a new class of psychological games. Section 4 presents the notion of psychological sequential equilibrium. Section 5 contains the interactive epistemology analysis. Section 6 discusses extensions which thus far have not been covered, including how to deal with incomplete information. Section 7 concludes. An appendix collects some of the proofs.

## 2 Overview of the conceptual issues

This section surveys the conceptual issues that motivate our work. After a preamble, we go through **R1**-**R4** (from the Introduction) in more detail, and provide supporting examples. The style is 'semi-technical'; we introduce some notation, but postpone a proper treatment of details for later sections.

The traditional approach to analyzing extensive games describes a player's preferences using a utility function of the form

---

[6]Kolpin (1992) explores an alternative route to analyzing GPS' games, in which players 'choose beliefs'. Gul & Pesendorfer (2004) propose an alternative framework to model social preferences that does not feature belief-dependent motivations, and yet is able to capture some of the phenomena typically modeled with psychological games. Segal & Sobel (2003) analyze simultanous moves games, and assume that preferences over material consequences depend on the equilibrium probability distribution over actions. They show that their approach can be regarded as a reformulation of GPS' normal form games.

$$u_i : Z \to \mathbb{R}$$

where $Z$ is the set of terminal histories (end nodes).

Psychological games are designed to capture richer motivations than traditional games, and the payoff functions have richer domains. GPS define a set of $i$'s initial (pre-play) beliefs about others' strategies and initial beliefs, here referred to as $\overline{\mathbf{M}}_i$, which does not rule out any hierarchy of initial beliefs. GPS model preferences using utility functions of the form

$$u_i : Z \times \overline{\mathbf{M}}_i \to \mathbb{R}$$

Their approach is rich enough to model interesting forms of belief-dependant motivation. Example 1 of the Introduction, $e.g.$, could be handled by assuming that Ann's payoff equals to $w - t - 2|\tau - t|$, where $w$ is her pre-tip wealth, $t \in \{0, 1, ..., w\}$ is her tip, and $\tau$ is her expectation of Bob's expectation of $t$. Ann would maximize her payoff by choosing $t = \tau$.[7]

However, the issues **R1**-**R4** lead us to enrich the domain of utilities further. We consider payoff functions of the form

$$u_i : Z \times \prod_{j \in N} \mathbf{M}_j \times \prod_{j \in N} S_j \to \mathbb{R}$$

where $\mathbf{M}_j$ is the set of $j$'s possible *conditional* beliefs about others' strategies and conditional beliefs, $S_j$ is the set of (pure) strategies of $j$, and $N$ is the set of players. The conditioning in $\mathbf{M}_j$ is done for every history, building on Battigalli & Siniscalchi (1999) who show how to represent hierarchies of conditional beliefs without ruling out any hierarchy. $\overline{\mathbf{M}}_j$ is (isomorphic to) a subspace of $\mathbf{M}_j$, so the payoff functions we consider are more general than those assumed by GPS.[8]

It is useful to keep these functional forms in mind as we go, because the issues **R1**-**R4** can be related to different arguments of $u_i$.

**R1 (updated beliefs):**

Rabin's (1993) theory of reciprocity, in which players reciprocate belief-dependent (un)kindness with (un)kindness, is probably the most well-known application of GPS' theory. Rabin works with the normal form version of GPS' theory. His goal is to highlight certain key qualitative features of reciprocity, and he does not address issues of dynamic decision making although

---

[7]See Charness & Dufwenberg (2004, ex. 1) for more discussion of a similar example.

[8]For a more precise comparison between our framework and GPS see subsection 6.1.

he points out that this is important for applied work (p. 1296). Dufwenberg & Kirchsteiger (2004) pick up from there, and develop a theory of reciprocity for extensive games. In motivating their exercise, they argue that it is necessary to deviate from GPS' extensive form framework: GPS only allow initial beliefs to enter the domain of a player's utility, while the modeling of reciprocal response at various ventures of a game tree requires that kindness be re-evaluated using updated belief. The argument is an instance of **R1**.

Reciprocity theory does not provide the easiest route to illustrating the key issues involved, however. Instead, we will consider the motivation of guilt aversion, applied to the trust game $\Gamma_1$.[9] Payoffs are in dollars and do not necessarily represent preferences. For this reason we call them 'material payoffs'.
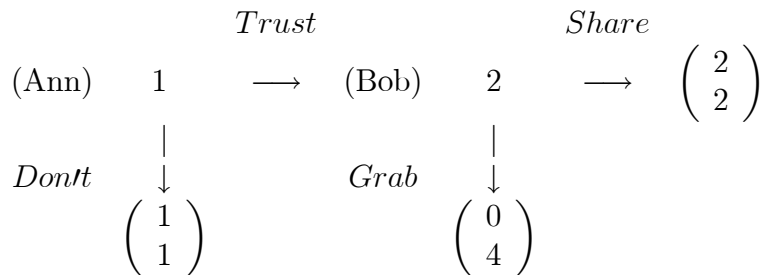
$$
\begin{array}{cccc}
& \textit{Trust} & & \textit{Share} \\
\text{(Ann)} \quad 1 & \longrightarrow & \text{(Bob)} \quad 2 & \longrightarrow \quad \begin{pmatrix} 2 \\ 2 \end{pmatrix} \\
\textit{Don!t} \quad \big\downarrow & & \textit{Grab} \quad \big\downarrow & \\
\begin{pmatrix} 1 \\ 1 \end{pmatrix} & & \begin{pmatrix} 0 \\ 4 \end{pmatrix} &
\end{array}
$$

**Figure 1.** The Trust Game $\Gamma_1$ with material payoffs

We now modify $\Gamma_1$ to incorporate a guilt sentiment of Bob's: Let $\alpha$ be the probability that Ann (initially) assigns to Bob's strategy *Share if Trust*. Bob suffers from guilt to the extent that he believes he lets Ann down. He argues that the higher is $\alpha$ the more let down she will be if he chooses *Grab*. Bob does not know what $\alpha$ is, as this belief is in the mind of Ann. However, he has a belief about $\alpha$. Let $\beta$ be Bob's expectation of $\alpha$, conditional on Ann choosing *Trust*. We can model guilt aversion assuming that Bob's utility at the terminal history (*Trust, Grab*) is decreasing in $\beta$.

The psychological game $\Gamma_2$ models this. What appears at the terminal histories should be thought of as utilities, not as material payoffs although the notions coincide for all but one terminal histories.[10]

---

[9]Charness & Dufwenberg's (2004), coin the term "guilt aversion" and develop a theory of it within the framework of GPS. Huang & Wu (1994), Dufwenberg (1995, 2002), Dufwenberg & Gneezy (2000), Bacharach, Guerra & Zizzo (2002), and Guerra & Zizzo (2004) consider related sentiments in trust games.

[10]There is no special significance to the "5" in Figure 2; we could have chosen many other numbers to get a working numerical example. Similar remarks apply to all the examples that follow in this section.

$$
\begin{array}{ccccccc}
 & Trust & & & & Share & \\
\text{(Ann)} \quad 1 & \longrightarrow & \text{(Bob)} & 2 & \longrightarrow & \begin{pmatrix} 2 \\ 2 \end{pmatrix} \\
\end{array}
$$

Don*t $\quad\downarrow \qquad\qquad\quad$ Grab $\quad\downarrow$

$$
\begin{pmatrix} 1 \\ 1 \end{pmatrix} \qquad\qquad\qquad \begin{pmatrix} 0 \\ 4-5\beta \end{pmatrix}
$$

**Figure 2.** The Psychological Trust Game $\Gamma_2$

$\Gamma_2$ is not a game in GPS' class, because $\beta$ (being an updated belief) is not captured by any element of $\overline{\mathbf{M}}_i$. This in itself illustrates **R1**. However, in order to appreciate the significance of this issue, it is useful to note that one can draw compelling (we think) conclusions about behavior that hinge crucially on the fact that $\beta$ is an updated belief.

Following Dufwenberg (1995, 2002), consider the following (for the time being intuitive) 'psychological forward induction' argument: Suppose Ann chooses *Trust*. If she is rational, this implies that she believes the probability that Bob would choose *Share* (after *Trust*) is at least $\frac{1}{2}$, *i.e.*, $\alpha \geq \frac{1}{2}$. Since we can figure this out, presumably Bob can too. Even if he is uncertain regarding the relevant value of $\alpha$, he infers it is at least $\frac{1}{2}$. Hence $\beta \geq 1/2$. Since $4 - 5\beta < 2$ if $\beta \geq \frac{1}{2}$, he prefers *Share*. Since we can figure this out, presumably Ann can too. Hence she chooses *Trust*, fully expecting Bob to *Share* (so $\alpha = 1$). Bob figures this out (so that $\beta = 1$), which further reinforces his preference to *Share*. Thus, the path (*Trust, Share*) is predicted!

The logic of the preceding argument depends on belief $\beta$ being *conditional* on Ann choosing *Trust*. The argument cannot be recast using GPS' theory, since $\overline{\mathbf{M}}_i$ contains only initial beliefs, but it can be captured in our framework, since $\mathbf{M}_i$ contains all hierarchies of conditional beliefs.

We close our discussion of **R1** by arguing that once one allows for conditional beliefs to enter utilities, it becomes natural to include beliefs conditional on *terminal* histories (which is allowed by our framework). Consider the following example $\Gamma_3$. We only specify utilities for player 1 (Ann), because the the point we wish to make is independent of player 2's (Bob) payoffs.

**Figure 3.** The Psychological Regret Game $\Gamma_3$

$\Gamma_3$ is a psychological game, in which Ann is disposed to suffer from regret, which we model via belief-dependent utility components at the terminal histories $(\ell, L)$ and $(r, L')$.[11] $\gamma$ is the probability with which Ann, conditional on $(\ell, L)$ being reached, thinks that Bob would have chosen $R'$ had she chosen $r$; $\delta$ is the probability with which Ann, conditional on $(r, L')$, thinks Bob would have chosen $R$ had she chosen $l$. The idea is that Ann's regret depends on what she thinks would have happened had she chosen differently, an assessment which might change in light of new information such as a choice by Bob.[12] The crucial thing to note in regard to **R1** is that $\gamma$ and $\delta$ are *conditional* beliefs at *terminal* histories.[13]

### R2 (others' beliefs):

There are two independent justifications for letting a player's utility depend on others' beliefs. First, this may be an adequate description of how certain social rewards operate. Refer back to example 3 from the Introduction, where Eva's preferences over effort depends on Fred's inferences. The example is (essentially) taken from Dufwenberg & Lundholm (2001). A related example is Bernheim's (1994) model of social conformity. These authors develop formal models where a player's utility depends directly on others' beliefs (although they do not make reference to psychological games).[14]

---

[11]The notation $(\ell, L)$ and $(r, L')$ refers to the sequence of action labels involved.

[12]Bell (1982) and Loomes & Sugden (1982) develop theories of regret, in which a decision maker's experienced utility depends on the post-choice revelation of a state-of-nature. Our formulation preserves that spirit, but extends it to belief-dependent motivation. This is natural in a strategic setting, where players cannot perfectly observe *ex post* the state of the world, which includes what another player *would* have chosen.

[13]In $\Gamma_3$, Ann's utility depends directly on her terminal history belief. By expanding on the theme, one can readily cook up examples where a player's utility depends on beliefs about a terminal history belief of another player. For example, Bob's utility at $(\ell, L)$ could depend on $\varepsilon$, defined as his expectation of $\delta$, conditional on history $r$.

[14]Indeed, these models can be interpreted as psychological games with asymmetric information where the utility of the informed player depends on the terminal beliefs of the uninformed player (see subsection 6.2).

The second justification concerns convenience in modeling. Refer back to the discussion concerning $\Gamma_2$, including the definition of $\alpha$ and $\beta$. Recall that in $\Gamma_2$ we modeled Bob's guilt feelings by letting his psychological payoff depend on $\beta$. It turns out that one has an equivalent modeling choice. Rather than go with $\Gamma_2$, one can assume that Bob's utility at $(\ell, L)$ depends directly on $\alpha$, rather than on $\beta$, although Bob is uncertain about the true value of $\alpha$ so that he uses probability assessments to weigh the different possibilities. We then get $\Gamma_4$:



$$
\begin{array}{ccccc}
 & \textit{Trust} & & \textit{Share} & \\
(\text{Ann}) \quad 1 & \longrightarrow & (\text{Bob}) \quad 2 & \longrightarrow & \begin{pmatrix} 2 \\ 2 \end{pmatrix} \\
\textit{Don{\textquotesingle}t} \quad \downarrow & & \textit{Grab} \quad \downarrow & & \\
\begin{pmatrix} 1 \\ 1 \end{pmatrix} & & \begin{pmatrix} 0 \\ 4 - 5\alpha \end{pmatrix} & &
\end{array}
$$

**Figure 4.** The Psychological Trust Game $\Gamma_4$

After *Trust*, when Bob has to make a choice he compares 2, the payoff of action *Share*, with the conditional expected payoff of action *Grab*, that is $E_2[4 - 5\alpha | Trust] = 4 - 5\beta$; we obtain the same results as with $\Gamma_2$.[15]

This example illustrates a general point: some belief-dependent motivations can be modeled replacing a conditional own belief of a certain 'order' (meaning: how many layers of beliefs about beliefs/choices are involved) with another object involving one degree lower order.[16] This may allow one to work with utilities of the form $u_i : Z \times \prod_{j \neq i} \mathbf{M}_j \times \prod_{j \in N} S_j \to \mathbb{R}$, where $\mathbf{M}_i$ is *not* a factor of the domain. This has two methodological advantages. First, it may seem easier to work with lower order beliefs (like having a first-order belief like $\alpha$ rather than a second-order belief like $\beta$). Second, and most importantly, one is lead to clearly distinguish between the carriers of utility (*i.e.*, elements of $Z \times \prod_{j \neq i} \mathbf{M}_j \times \prod_{j \in N} S_j$) and how a player deals with uncertainty by making probabilistic predictions and updating them (described by elements of $\mathbf{M}_i$). By contrast, when the domain of $i$'s utility is $Z \times \prod_{j \in N} \mathbf{M}_j \times \prod_{j \in N} S_j$ elements of $\mathbf{M}_i$ end up serving both purposes.

---

[15] We do not suggest that $\Gamma_4$ is interesting only in that it provides a convenient alternative way to analyzing $\Gamma_2$; the emotion modeled in $\Gamma_4$ may make sense in its own right, as a primitive assumption about preferences (akin to example 3 of the introduction).

[16] A special case of this appears if a player holds a conditional belief about a choice, in which case the object involving one degree lower order would be that choice itself. The object may also be an initial belief, like $\alpha$ in the example.

## R3 (dependence on strategies):

Many forms of belief-dependent motivation are compelling in combination with preferences that depend on overall strategies, beyond what is captured by how strategies cause terminal histories to be reached. Consider the psychological game $\Gamma_5$, a variation of $\Gamma_1$ where Ann may 'dissipate' all the material payoffs. The payoffs of $\Gamma_5$ are material, not necessarily reflecting utilities.



**Figure 5.** Modified Trust Game $\Gamma_5$ with Material Payoffs

To make our point, let us first model a sentiment of Ann's: the stronger she expects Bob to *Share*, the more let down she feels at the terminal history (*Trust, Grab*), and this feeling is painful to her. One way to model this is to assume that her utility at (*Trust, Grab*) is not 0 (as in $\Gamma_5$), but rather $0 - 5\alpha = -5\alpha$, where $\alpha$ is the probability Ann assigns to Bob's strategy *Share if Trust*. However, arguably this assumption has the following flaw: Suppose that Ann is planning to choose *Dissipate*. In this case, she is bound to get zero material payoff whether or not Bob chooses *Share*, so it may makes little sense for her to feel disappointed if he does! We propose that a natural reaction is to let Kate's utility following (*Trust, Grab*) be $-5\alpha$ if she plans to choose *Keep*, but 0 if she plans to choose *Dissipate*. In this case, Ann's utility depends *on her own choice of strategy*, on top of which terminal history is reached and which beliefs she harbors.

A slight further complication of the example will suggest that a player's utility may also reasonably depend on another player's strategy. To this end, focus on Bob. Assume that he suffers from guilt to the extent that he believes he lets Ann down, so that he suffers a utility loss following (*Trust,Grab*). One way to model this is to let his utility at (*Trust, Grab*) be $4 - 5\alpha$ (as in $\Gamma_4$). This specification would, however, seem to suffer from a flaw analogous to the one we discussed in the previous paragraph.[17] A natural reaction is to

---

[17]Suppose that Gwen is planning to choose *Dissipate*. In that case, she is bound to get no material payoff whether or not Hugh chooses *Share*, so it makes little sense for him to feel that he lets her down if he does.

modify Bob's payoff following (*Trust,Grab*) to be $4 - 5\alpha$ if Ann plans to choose *Keep*, but 4 if she plans to choose *Dissipate*.

These examples to illustrate the following: Psychological motivations often exhibit a concern, not only for players' actual actions, but also for their intentions. Intentions depend on beliefs as well as on strategies, and the latter dependence goes beyond what is implied by how strategies induce endnodes. Therefore, the domain of our psychological utility function includes (conditional) beliefs and strategies of every player, on top of terminal nodes.

## R4 (non-equilibrium analysis):

**R1**-**R3** concern features of players' motivation that one may wish to incorporate in a formal framework. The next step is to generate predictions about behavior. We propose a notion of psychological sequential equilibrium (PSE), which generalizes the sequential equilibrium concept of Kreps & Wilson (1982). We postpone illustrations of PSE until we have formally introduced the concept in section 4.

While much of economic theory presumes that players coordinate on an equilibrium, it is not always clear that such an assumption is justified. For one thing, people may be quite rational, and confident in others' rationality, even if they fail to coordinate on an equilibrium. In conventional game theory, related matters have inspired work on the implications of common belief of rationality; see *e.g.* the work by Bernheim (1984) and Pearce (1984) on rationalizability. This brings us to **R4**. There is little reason to assume that equilibrium coordination is easier in psychological games than in standard games. In fact, since psychological games often seem more complicated, and since problems of equilibrium multiplicity are likely to be enhanced in psychological games, assuming equilibrium may be assuming too much *especially* in psychological games.

Another reason to feel skeptical about a fully fledged equilibrium analysis in psychological games is the following: It is often argued that players learn to play Nash equilibrium because through recurrent strategic interaction they come to hold correct beliefs about the actions of the opponents (see, *e.g.*, Fudenberg & Levine, 1998, and references therein). This is *not enough* for a psychological equilibrium; since payoffs depend on hierarchical beliefs, players would have to be able to learn the beliefs of others, but unlike actions beliefs are typically not observable *ex post*.

Giving up the equilibrium assumption does not, however, necessarily mean giving up on predictive power. Refer back to the psychological forward induction argument, presented in conjunction with $\Gamma_2$. Ann and Bob were presented as performing deductive reasoning regarding one another's

behavior and beliefs, and a clear-cut prediction resulted despite that no presumption of equilibrium was made. However, the story told was informal, and specific to $\Gamma_2$ (or, equivalently, $\Gamma_4$). It is natural to wonder about formalizations that are generally applicable. In section 5, we develop a framework for analyzing interactive epistemology in psychological games, without postulating equilibrium play. In particular, building on an epistemic theme due to Battigalli & Siniscalchi (2002), we extend Pearce's (1984) classical notion of (extensive form) rationalizability to psychological games. The concept captures psychological forward induction in simple games like $\Gamma_2$ and $\Gamma_4$, and in more complicated games for which long chains of beliefs about beliefs may be needed to get clear-cut predictions.

# 3   Psychological Games

In this section we develop a formal framework to analyze dynamic psychological games. We introduce the notation on extensive-form games (3.1), model a universal belief space that accounts for updating beliefs about others' beliefs (3.2), and put forth and illustrate our general definition of a psychological game (3.3).

## 3.1   Extensive forms with observable actions

To simplify the analysis we restrict our attention to finite multi-stage games with observable actions and no chance moves. For the time being, we also rule out incomplete information. These restrictions can be removed, at the cost of additional complexity in notation (see section 6). We assume that players move simultaneously at every stage of the game. This is without loss of generality, because the set of feasible actions of a player may depend on the actions chosen in previous stages and it may be a singleton. Simultaneous moves games, perfect information games and repeated games are special cases (cf. Osborne & Rubinstein, 1994, ch. 6). We use the following notation and terminology:[18]

An *extensive form* with observable actions is a tuple $\langle N, H \rangle$ where $N = \{1, ..., n\}$ is the *player* set, and $H$ is the set of feasible *histories* of the game. A history of length $\ell$ is a sequence $h = (a^1, ..., a^\ell)$ where each $a^t = (a_1^t, ..., a_n^t)$ represents the profile of actions chosen at stage $t$ ($1 \leq t \leq \ell$). We assume that a history $h$ becomes public information as soon as it occurs. We also assume that $H$ is finite. For notational convenience, we let $H$ contain the *empty history*, denoted by $h^0$ (the history of length 0). The set of feasible

---

[18]See the Appendix for a more rigorous and complete definition of each term.

actions for player $i$ at history $h$ is denoted by $A_i(h)$ and it may be a singleton, meaning that $i$ is not active at $h$. $A_i(h)$ is empty if and only if $h$ is a *terminal* history. We let $Z$ denote the set of terminal histories.

Extensive forms with only one active player at each nonterminal history are graphically represented by trees, following standard conventions (see the examples in section 2).

For any given extensive form, we let $S_i$ denote the set of (pure) strategies of player $i$. A typical strategy is denoted by $s_i = (s_{i,h})_{h \in H \setminus Z}$, where $s_{i,h}$ is the action that would be selected by strategy $s_i$ if history $h$ occurred. Define $S = \prod_{i \in N} S_i$ and $S_{-i} = \prod_{j \neq i} S_j$. The set of strategies of player $i$ that allow history $h$ is denoted $S_i(h)$. A similar notation is used for strategy profiles: $S(h) = \prod_{i \in N} S_i(h)$ and $S_{-i}(h) = \prod_{j \in N} S_j(h)$. Finally, we let $\zeta(s) \in Z$ denote the terminal history induced by strategy profile $s = (s_i)_{i \in N}$.

## 3.2 Conditional beliefs & infinite hierarchies of beliefs

Here we summarize the theory of hierarchies of conditional beliefs due to Battigalli & Siniscalchi (1999), which should be consulted for proofs, details and further references. Consider a decision maker DM who is uncertain about which element in a set $X$ is true. Assume $X$ is a compact Polish space.[19] DM assigns probabilities to events $E$, $F$, ... in the Borel sigma-algebra $\mathcal{B}$ of $X$ according to some (countably additive) probability measure. Let $\Delta(X)$ denote the set of all such probability measures. As events unfold DM updates her beliefs in a coherent fashion. The actual and/or potential beliefs of DM are described by a conditional probability system (see Rênyi, 1955). Let $\mathcal{C} \subseteq \mathcal{B}$ denote the collection of potentially observable events (or conditioning events). DM holds probabilistic beliefs conditional on each event $F \in \mathcal{C}$.

**Definition 1** *A* conditional probability system *(cps) on* $(X, \mathcal{B}, \mathcal{C})$ *is a function* $\mu(\cdot|\cdot) : \mathcal{B} \times \mathcal{C} \to [0, 1]$ *such that for all* $E \in \mathcal{B}$, $F, F' \in \mathcal{C}$
*(1)* $\mu(\cdot|F) \in \Delta(X)$,
*(2)* $\mu(F|F) = 1$,
*(3)* $E \subseteq F' \subseteq F$ *implies* $\mu(E|F) = \mu(E|F')\mu(F'|F)$.

We regard the set of cps' on $(X, \mathcal{B}, \mathcal{C})$ as a subset of the topological space $[\Delta(X)]^{\mathcal{C}}$, where $\Delta(X)$ is endowed with the topology of weak convergence of measures and $[\Delta(X)]^{\mathcal{C}}$ is endowed with the product topology.

¿From now on DM is a player $i$, and $(X, \mathcal{B}, \mathcal{C})$ is specified as follows: either $X = S_{-i}$ (a finite set), or $X = S_{-i} \times Y$, where $Y$ is some compact Polish

---

[19]A topological space $X$ is *Polish* if it admits a compatible metric $d$ such that $(X, d)$ is a complete and separable metric space (see, *e.g.*, Kechris, 1995, p 13).

parameter space typically representing a set of opponents' beliefs; the Borel sigma-algebra $\mathcal{B}$ is implicitly understood,[20] and conditioning events corresponds to histories, that is, $\mathcal{C} = \{F \subseteq S_{-i} \times Y : F = S_{-i}(h) \times Y, h \in H\}$ (or $\mathcal{C} = \{F \subseteq S_{-i} : F = S_{-i}(h), h \in H\}$ if $X = S_{-i}$). Accordingly, the set of cps' is denoted $\Delta^H(S_{-i} \times Y)$ a subset of $[\Delta(S_{-i} \times Y)]^H$. We use the following abbreviation: if conditioning event $F$ corresponds to history $h$ then we write $\mu(\cdot|F) = \mu(\cdot|h)$.

Now take the point of view of an opponent of player $i$, who is uncertain about the true (strategy and) cps of other players. The following result shows that we can take for granted that $\Delta^H(S_{-i} \times Y)$ is a compact Polish space, like the given parameter space $Y$.[21] This result is key in our construction of hierarchical conditional beliefs in that it shows that the domain of higher-order uncertainty has the same structural properties as the domain of lower-order uncertainty.

**Lemma 2** $\Delta^H(S_{-i})$ *is a compact Polish space. Furthermore, if $Y$ is a compact Polish space, also $\Delta^H(S_{-i} \times Y)$ is a compact Polish space.*

Hierarchies of cps' are defined recursively as follows:

- $X_{-i}^0 = S_{-i}$ $(i \in N)$,

- $X_{-i}^k = X_{-i}^{k-1} \times \prod_{j \neq i} \Delta^H(X_{-j}^{k-1})$ $(i \in N; k = 1, 2, ...)$.

By repeated applications of Lemma 2, each $X_{-i}^k$ is a cross-product of compact Polish spaces, and hence it is a compact Polish space in itself.[22] A cps $\mu_i^k \in \Delta^H(X_{-i}^{k-1})$ is called $k$-order cps. For $k > 1$, $\mu_i^k$ is a joint cps on the strategies and $(k-1)$-order cps' of the opponents. A *hierarchy of cps'* is a countably infinite sequence of cps' $\boldsymbol{\mu}_i = (\mu_i^1, \mu_i^2, ...) \in \prod_{k>0} \Delta^H(X_{-i}^{k-1})$. Hierarchy $\boldsymbol{\mu}_i$ is *coherent* if the cps' of distinct orders assign the same conditional probabilities to lower-order events, that is

$$\mu_i^k(\cdot|h) = \mathrm{marg}_{X_{-i}^{k-1}} \mu_i^{k+1}(\cdot|h) \ (k = 1, 2, ...; \ h \in H).$$

It can be shown that a coherent hierarchy $\boldsymbol{\mu}_i$ induces a cps $\nu_i$ on the cross-product of $S_{-i}$ with the sets of hierarchies of beliefs of $i$'s opponents, a compact Polish space.

---

[20] $\mathcal{B}$ is obtained from the product of the discrete topology on $S_{-i}$ and the topology of $Y$.

[21] This depends on two facts: (1) the collection of conditioning events for player $i$ (corresponding to $H$) is at most countable (indeed finite), and (2) each conditioning event $S_{-i}(h) \times Y$ (or $S_{-i}(h)$ if $X = S_{-i}$) is both closed and open.

[22] The cross-product of countably many compact Polish spaces is also compact Polish.

However, $\nu_i$ may assign positive probability (conditional on some $h$) to opponents' incoherence. To rule this out, say that a coherent hierarchy $\boldsymbol{\mu}_i$ satisfies belief in coherency if the induced cps $\nu_i$ is such that each $\nu_i(\cdot|h)$ ($h \in H$) assigns probability one the opponents' coherency; $\boldsymbol{\mu}_i$ satisfies belief in coherency of order $k$ if it satisfies belief in coherency of order $k-1$ and the induced cps $\nu_i$ is such that each $\nu_i(\cdot|h)$ ($h \in H$) assigns probability one the opponents' coherency of order $k-1$; $\boldsymbol{\mu}_i$ is *collectively coherent* if it satisfies belief in coherency of order $k$ for each positive integer $k$. The set of collectively coherent hierarchies of player $i$ is a compact Polish space, denoted by $\mathbf{M}_i$. We let $M_i^k$ denote the set of of $k$-order beliefs consistent with collective coherency, that is, the projection of $\mathbf{M}_i$ on $\Delta^H(X_{-i}^{k-1})$, and let $M_{-i}^k = \prod_{j\neq i} M_j^k$, $\mathbf{M}_{-i} = \prod_{j\neq i} \mathbf{M}_j$, $\mathbf{M} = \prod_{j\in N} \mathbf{M}_j$.

We have now defined all the elements that form the domain of the psychological utility functions. But is this enough for the analysis of strategic reasoning? In order to decide on the best course of action player $i$ may need to form (conditional) beliefs about the *infinite* hierarchies of (conditional) beliefs of other players, either because they enter his psychological payoff function or because his assessment of the behavior and finite-order beliefs of other players is derived from assumptions, such as "common belief in rationality", involving beliefs of infinitely many orders. Does this mean that we need additional layers of beliefs? No. The following result shows that the countably infinite hierarchies of cps' defined above are sufficient for the strategic analysis because $\mathbf{M}_i$ is isomorphic to $\Delta^H(S_{-i} \times \mathbf{M}_{-i})$, hence each $\boldsymbol{\mu}_i \in \mathbf{M}_i$ corresponds to a cps on $S_{-i} \times \mathbf{M}_{-i}$:

**Lemma 3** *For each player $i \in N$ there is a one-to-one and onto continuous function*

$$f_i = (f_{i,h})_{h\in H} : \mathbf{M}_i \to \Delta^H(S_{-i} \times \mathbf{M}_{-i})$$

*whose inverse is also continuous. Furthermore, each coordinate function $f_{i,h}$ is such that for all $\boldsymbol{\mu}_i = (\mu_i^1, \mu_i^2...) \in \mathbf{M}_i$, $k \geq 1$*

$$\mu_i^k(\cdot|h) = \mathrm{marg}_{S_{-i} \times M_{-i}^1 \times ... \times M_{-i}^{k-1}} f_{i,h}(\boldsymbol{\mu}_i).$$

## 3.3   Psychological Games

We are now ready to state our definition of a psychological game:

**Definition 4** *A* psychological game *based on extensive form $\langle N, H \rangle$ is a structure $\Gamma = \langle N, H, (u_i)_{i\in N} \rangle$ where $u_i : Z \times \mathbf{M} \times S \to \mathbb{R}$ is $i$'s (measurable and bounded) psychological payoff function.*

The numerical examples examined in section 2 fit this definition: in game $\Gamma_2$, $u_2$ depends on $z$ and $\mu_2^2(\cdot|Trust)$;[23] in game $\Gamma_3$, $u_1$ depends on $z$ and $\mu_1^1(\cdot|z)$; in game $\Gamma_4$, $u_2$ depends on $z$ and the initial first-order belief of player 1, $\mu_1^1(\cdot|h^0)$; finally, the psychological payoff functions related to $\Gamma_5$ (a game with material payoffs) let $u_1$ and $u_2$ depend on $z$, $\mu_1^1(\cdot|h^0)$, and $s_1$.

In all these examples, a psychological game is obtained from a *material payoff game* $\langle N, H, (\pi_i : Z \to \mathbb{R})_{i \in N} \rangle$ according to some formula that captures psychological motivations like regret, feeling let down, or guilt. To illustrate our framework we provide a few instances of such derivations. We focus on two-person games.

We can obtain psychological game $\Gamma_3$ by adding a regret component to the material payoff of player 1 (Ann). Regret of player $i$ at a terminal history $z$ can be captured by the difference between the actual material payoff $\pi_i(z)$ and the maximal expected payoff that could have been obtained 'with the benefit of hindsight,' that is, using the terminal beliefs conditional on $z$. Formally, we can measure regret with the (negative) function

$$R_i(z, \boldsymbol{\mu}_i^1(\cdot|z)) = \pi_i(z) - \max_{s_i} \sum_{s_j' \in S_j(z)} \mu_i^1(s_j'|z)\pi_i(\zeta(s_i, s_j'))$$

and we obtain the psychological payoff function

$$u_i(z, \boldsymbol{\mu}, s) = \pi_i(z) + \theta_i R_i(z, \boldsymbol{\mu}_i^1(\cdot|z))$$

where $\theta_i$ is a psychological sensitivity parameter. This formulation can then be applied to incorporate regret to any extensive game with a given material payoff function. In game $\Gamma_3$, *e.g.*, we assumed that $\pi_1(\ell, L) = \pi_1(r, L') = 0$, $\pi_1(\ell, R) = \pi_1(r, R') = 1$, and $\theta_1 = 1$.

As exemplified with reference to the material payoff game $\Gamma_5$, player $i$ feels 'let down' if her actual material payoff $\pi_i(z)$ is lower than the payoff she expected to get, given her initial first-order beliefs $\mu_i^1(\cdot|h^0)$ and her strategy $s_i$. This can be measured by the negative function

$$D_i(z, \mu_i^1(\cdot|h^0), s_i) = \min\left\{0, [\pi_i(z) - \sum_{s_j'} \mu_i^1(s_j'|h^0)\pi_i(\zeta(s_i, s_j'))]\right\}.$$

Aversion to this feeling can be captured by the psychological payoff function

$$u_i(z, \boldsymbol{\mu}, s) = \pi_i(z) + \theta_i D_i(z, \mu_i^1(\cdot|h^0), s_i),$$

---

[23]$\mu_2^2(\cdot|Trust)$ is the conditional second-order belief of player 2 (Bob) used to compute the expectation $\beta$ of $\alpha$, the probability initially assigned by player 1 (Ann) to the strategy *'Share if Trust'*.

and a guilt motivation can be modeled as aversion to letting the other player down, which can be captured by the following payoff function:[24]

$$u_i(z, \boldsymbol{\mu}, s) = \pi_i(z) + \theta_i D_j(z, \mu_j^1(\cdot|h^0), s_j)$$

For the special case of the Trust Game, we obtain psychological game $\Gamma_4$ by letting $\theta_1 = 0$ and $\theta_2 = \frac{5}{2}$.

# 4 Equilibrium Analysis

Kreps & Wilson's (1982) notion of sequential equilibrium has become a benchmark for the analysis of standard extensive games. Our goal here is to extend this concept to the class of psychological games defined in section 3. (The restriction to multi-stage games with complete information simplifies the presentation, but is actually not essential as we discuss in section 6.) We next comment on the entailed interpretation of mixed strategies and assessments (4.1), give the main definition (4.2), consider some examples (4.3).

## 4.1 Randomized strategies and consistent assessments

The equilibrium concept we develop refers to randomized choices. However, in our interpretation, we exclude explicit randomization (players tossing coins or spinning roulette wheels). Rather, we will interpret a randomized choice of a given player $i$ as the common first-order belief of $i$'s opponents about $i$ (cf. Aumann & Brandenburger, 1995). This is the analog of the following characterization of a Nash equilibrium in a standard simultaneous moves game: a profile $(\sigma_1, ..., \sigma_n) \in \Delta(A_1) \times ... \times \Delta(A_n)$ is an equilibrium if for each player $i$ each action in the support of $\sigma_i$ is a best response to $\sigma_{-i}$.

In the analysis of extensive form games we focus on behavior strategies (rather than mixed strategies): $\sigma_i = (\sigma_i(\cdot|h))_{h \in H \setminus Z} \in \prod_{h \in H \setminus Z} \Delta(A_i(h))$. We interpret a behavior strategy $\sigma_i$ as an array of common conditional first-order beliefs held by the opponents of player $i$. This interpretation is part of the notion of 'consistency' of profiles of behavior strategies and hierarchical beliefs defined below.

---

[24]We have opted here for a formulation such that $i$'s psychological payoff depends directly on others' beliefs; cf. the discussion of **R2** and $\Gamma_2$ vs. $\Gamma_4$ in Section 2. Moreover, only initial beliefs enter the utility function directly. However, in the strategic analysis the updated second-order beliefs of $i$ are crucial because they determine the expected payoff that $i$ maximizes at each history. The reader may want to compare the formulation here to the related one by Charness & Dufwenberg (2004), which by contrast utilizes GPS' framework and hence lets only $i$'s own initial beliefs influence $i$'s utility.

Kreps & Wilson (1982) argue that an appropriate definition of equilibrium in extensive form games must refer to 'assessments,' that is, profiles of (behavior) strategies *and* conditional (first-order) beliefs. They formulate a definition of sequential equilibrium in two steps: first they put forward a 'consistency' condition for assessments, and then they stipulate that an assessment is a sequential equilibrium if it is consistent and it satisfies sequential rationality. It turns out that the consistency condition captures the assumptions that ($a$) each player regards his opponents' strategies as stochastically independent, and ($b$) any two players have the same (prior and conditional) beliefs about any third player (cf. Fudenberg & Tirole 1991b, Battigalli 1996, and Kohlberg & Reny 1997). We follow a similar two-step approach, adding to it a third requirement concerning the higher-order beliefs that need to be specified in psychological games.

In our setup, an *assessment* is a profile $(\sigma, \boldsymbol{\mu}) = (\sigma_i, \boldsymbol{\mu}_i)_{i \in N}$ of behavior strategies and hierarchies of conditional beliefs. Before defining consistency of an assessments, we need to define more precisely what we mean by 'stochastic independence'. For this, we need to explain that a *marginal cps* on the strategies of $j$ is a cps on $(S_j, \mathcal{B}_j, \mathcal{C}_j)$, where $\mathcal{B}_j$ is the power set of $S_j$ and $\mathcal{C}_j = \{S_j(h), h \in H\}$. The set of such marginal cps' is denoted by $\Delta^H(S_j)$. The following definition takes advantage of the simple information structure we are assuming, *i.e.* perfect observability of past actions, and allows us to characterize stochastic independence for cps's in terms of 'marginal' cps's.

**Definition 5** *A first-order cps* $\mu_i \in \Delta^H(S_{-i})$ *satisfies the* stochastic independence *property, if there exists a profile of marginal cps's* $(\mu_{ij})_{j \neq i} \in \prod_{j \neq i} \Delta^H(S_j)$ *such that* $\mu_i(s_{-i}|h) = \prod_{j \neq i} \mu_{ij}(s_j|h)$ *for all* $h \in H$, $s_{-i} \in S_{-i}(h)$. *We let* $\Delta_I^H(S_{-i})$ *denote the set of first-order cps's of player* $i$ *that satisfy the stochastic independence property.*

Note that for each $\mu \in \Delta_I^H(S_{-i})$ we can derive a behavioral profile $(\sigma_j)_{j \neq i}$ as follows: let $S_j(h, a_j) = \{s_j \in S_j(h) : s_{j,h} = a_j\}$ denote the set of strategies of player $j$ that allow history $h$ and select action $a_j$ at $h$, then

$$\forall j \neq i, \forall h \in H, \forall a_j \in A_j(h), \ \sigma_j(a_j|h) = \mu_{ij}(S_j(h, a_j)|h). \qquad (1)$$

We are now ready for the main definition of this section:

**Definition 6** *A profile of hierarchies of cps'* $\boldsymbol{\mu} = (\boldsymbol{\mu}_i)_{i \in N} \in \mathbf{M}$ *is consistent if*
*(a) the first-order cps of each player satisfies stochastic independence, that is,*

$$\forall i \in N, \ \mu_i^1 \in \Delta_I^H(S_{-i}),$$

*(b) the marginal first-order beliefs of any two players about a third player coincide, that is*

$$\forall i, \forall j, \forall k \in N, \ \forall h \in H, \ marg_{S_k}\mu_i^1(\cdot|h) = marg_{S_k}\mu_j^1(\cdot|h),$$

*(c) each player's higher-order beliefs in $\boldsymbol{\mu}$ assign probability one to the lower-order beliefs in $\boldsymbol{\mu}$ itself, that is*

$$\forall i \in N, \ \forall k > 1, \ \forall h \in H, \ \mu_i^k(\cdot|h) = \mu_i^{k-1}(\cdot|h) \times \delta_{\mu_{-i}^{k-1}}$$

*where $\delta_x$ is the measure that assigns probability one to the singleton $\{x\}$. An assessment $(\sigma, \boldsymbol{\mu})$ is consistent if $\boldsymbol{\mu}$ is consistent and $\sigma$ is derived from first-order beliefs $(\mu_i^1)_{i \in N}$ as in eq. (1).*

The justification of the (very strong) conditions $(b)$ and $(c)$ comes from the classical interpretation of equilibrium beliefs: such beliefs are supposed to be the end-product of a transparent reasoning process that intelligent players can perform. Therefore any two players must share the same first-order conditional beliefs about any other player, and every player comes to a correct conclusion about the (hierarchical) beliefs of his opponents because he is able to replicate their reasoning process.[25] Condition $(c)$ is analogous to a condition used by GPS to define psychological Nash equilibrium. Essentially, it requires that players hold common, correct beliefs about each others' beliefs. This condition is equivalent to the requirement that, for each player $i$ and each history $h$, the conditional belief on $S_{-i} \times \mathbf{M}_{-i}$ induced by hierarchy $\boldsymbol{\mu}_i$ assigns probability one to $\boldsymbol{\mu}_{-i}$.[26]

## 4.2 Sequential Equilibrium Assessments

We take the point of view of an 'agent' $(i, h)$ of player $i$ who is in charge of the move at history $h$ and seeks to maximize $i$'s conditional expected payoff

---

[25]Condition $(c)$ yields the implication that, although players update their beliefs about the opponents' strategies, they never change their beliefs about what the opponents would believe conditional on each history. Of course, by observing the actual play-path each player infers the current actual beliefs of his opponents, but interesting forms of learning about the beliefs of others are ruled out. For example, condition $(c)$ implies that no player would ever change his mind about the initial beliefs of his opponents. Without defending this assumption, we argue that it is in the spirit of the standard definition of sequential equilibrium. After a hypothetical deviation by player $i$, this player is assumed to play a continuation strategy that maximizes his expected payoff against the same (equilibrium) beliefs that were ascribed to him before the deviation, even if the deviation is irrational under such beliefs (cf. Reny, 1992).

[26]That is, (3) holds iff $\forall i \in N, \ \forall h \in H, \ f_{i,h}(\boldsymbol{\mu}_i)(S_{-i} \times \{\boldsymbol{\mu}_{-i}\}) = 1$.

given the consistent belief profile $\boldsymbol{\mu}$. The expected payoff of $i$ conditional on history $h$ and action $a_i \in A_i(h)$ given $\mu$ can be expressed as

$$\mathrm{E}\ [u_i|h, a_i] = \sum_{s_{-i} \in S_{-i}(h)} \mu_i^1(s_{-i}|h) \sum_{s_i \in S_i(h, a_i)} \mu_{ji}^1(s_i|(h, a_i, s_{-i,h})) u_i(\zeta(s), \boldsymbol{\mu}, s)$$

$$(2)$$

where $s_{-i,h}$ is the action profile chosen by $i$'s opponents at $h$ according to $s_{-i}$ and $\mu_{ji}^1$ is the cps about $i$ of an arbitrary opponent $j$. This specification presumes that $(i, h)$ assesses the probabilities of actions by other agents of player $i$ in the same way as each player $j \neq i$; that explains how $\mu_{ji}^1(s_i|(h, a_i, s_{-i,h}))$ shows up in the right-hand-side of the expression.

The expected utility formula (2) is quite different from those used in the literature on standard games. This is because of the possibility that psychological payoffs are directly affected by strategies. When psychological payoffs are not directly affected by strategies $\mathrm{E}\ [u_i|h, a_i]$ can be expressed in a more familiar form:

**Remark 7** *Suppose that psychological payoff functions depend only on terminal histories and beliefs. Then for any consistent assessment $(\sigma, \boldsymbol{\mu})$*

$$\mathrm{E}\ [u_i|h, a_i] = \sum_z \Pr_\sigma[z|h, a_i] u_i(z, \boldsymbol{\mu})$$

*where $\Pr_\sigma[z|h, a_i]$ is the probability of terminal history $z$ conditional on $(h, a_i)$ determined by behavioral profile $\sigma$.*

We now move to the main definition of this section. A consistent assessment is a sequential equilibrium if it satisfies a sequential rationality condition:

**Definition 8** *Assessment $(\sigma, \boldsymbol{\mu}) = (\sigma_i, \boldsymbol{\mu}_i)_{i \in N}$ is a psychological sequential equilibrium (PSE) if it is consistent and for all $i \in N$, $h \in H$,*

$$\mathrm{Supp}(\sigma_i(\cdot|h)) \subseteq \arg \max_{a_i \in A_i(h)} \mathrm{E}\ [u_i|h, a_i]. \quad (3)$$

The sequential rationality condition (3) only requires that the assessment be immune to one-shot deviations. In standard games, and more generally in psychological games where payoff functions do not depend on own strategy, so that they have the form $u_i : Z \times \mathbf{M} \times S_{-i} \to \mathbb{R}$, the 'one-shot-deviation principle' applies. Condition (3) is then equivalent to requiring that the candidate equilibrium be immune to deviations to arbitrary continuation

strategies.[27] In subsection 5.3 we show that when the psychological payoff $u_i$ directly depends on $s_i$ the one-shot-deviation principle does not apply.

The main result of this section is an existence theorem:

**Theorem 9** *If the psychological payoff functions are continuous, there exists at least one psychological sequential equilibrium assessment.*

A complete proof is contained in the Appendix. Here we only provide a sketch. Existence can be shown by using 'Selten's trick' (Selten, 1975). Consider $\varepsilon$-perturbed games where there is positive minimal probability of choosing any action at any history, i.e. $\varepsilon = (\varepsilon_{i,h}(a_i, h)_{a_i \in A_i(h)})_{i \in N, h \in H}$ is a strictly positive vector such that $\sum_{a_i \in A_i(h)} \varepsilon(a_i, h) < 1$ for each history $h$. For each strictly positive behavior strategy profile, there exists a corresponding profile of hierarchies of cps' $\boldsymbol{\mu} = \beta(\sigma)$ such that $(\sigma, \beta(\sigma))$ is consistent.[28] For any $\varepsilon$-perturbed game we define an (agent-form, psychological) $\varepsilon$-equilibrium as an $\varepsilon$-constrained behavior strategy profile $\sigma_\varepsilon$ such that for each history $h$ and each player $i$, a pure action $a_i$ that does not maximize the expectation of $U_i$ (given $h$, $\beta(\sigma_\varepsilon)$ and $\sigma_\varepsilon$) is assigned the minimal probability $\varepsilon(a_i, h)$. It can be shown by standard compactness-continuity arguments that each $\varepsilon$-perturbed game has an $\varepsilon$-equilibrium (cf. the proof of existence of psychological Nash equilibria in GPS). Fix a sequence $\varepsilon^k \to 0$ and a corresponding sequence of $\varepsilon^k$-equilibrium assessments. By compactness, $\sigma^k$ has an accumulation point $\sigma^*$. By upper-hemicontinuity of the local best response correspondences, for each $(i, h)$, $\sigma_i^*(\cdot | h)$ assigns positive probability only to actions that are best responses to $(\sigma^*, \beta(\sigma^*))$ at $h$. Therefore $(\sigma^*, \beta(\sigma^*))$ is a psychological sequential equilibrium assessment.

We next show that the PSE concept generalizes the subgame perfect equilibrium concept of standard games with observable actions (recall that sequential and subgame perfect equilibrium coincide in games with observable actions). This is a consequence of a more general result about PSE in games where psychological payoffs depend only on terminal nodes and beliefs. Suppose that payoff functions depend only on the terminal history

---

[27]The 'one-shot-deviation principle' is essentially a dynamic programming result. It holds for finite (standard) games, and more generally for finite-horizon games, and infinite-horizon games where payoffs are 'continuous at infinity'. See, *e.g.*, Fudenberg & Tirole (1991a), pp 108-110.

[28]By Kuhn's transformation, a strictly positive behavior strategy profile $\sigma$ corresponds to a product measure $\mu_1 \times ... \times \mu_n$ on $S_1 \times ... \times S_n$; each marginal measure $\mu_{-i}$ on $S_{-i}$ yields a first-order cps for $i$ satisfying stochastic independence, and by construction the first-order cps's of different players agree; a corresponding profile of hierarchies is obtained assuming that there is 'common knowledge' of beliefs, as in condition ($c$) of Definition 6.

and hierarchical beliefs: $u_i : Z \times \mathbf{M} \to \mathbb{R}$. Then, for any fixed profile of hierarchies of cps' $\boldsymbol{\mu} = (\boldsymbol{\mu}_i)_{i \in N}$, we can obtain from a psychological game $\Gamma = \langle N, H, (u_i)_{i \in N} \rangle$ a standard game $\Gamma = \langle N, H, (v_i)_{i \in N} \rangle$ with payoff functions defined by $v_i(z) = u_i(z, \boldsymbol{\mu})$.

**Proposition 10** *Suppose that psychological payoff functions have the form $u_i : Z \times \mathbf{M} \to \mathbb{R}$. Then an assessment $(\sigma, \boldsymbol{\mu})$ is a psychological sequential equilibrium if and only if it is consistent and $\sigma$ is a subgame perfect (hence sequential) equilibrium of the standard game $\Gamma$.*

    **Proof.** First recall that when payoffs do not directly depend on strategies, the conditional expected payoffs determined by a consistent assessment $(\sigma, \boldsymbol{\mu})$ can be expressed as $\mathrm{E}[u_i | h, a_i] = \sum_z \Pr_\sigma[z | h, a_i] u_i(z, \boldsymbol{\mu})$ (see remark 7). Let $(\sigma, \boldsymbol{\mu})$ be a psychological sequential equilibrium. By definition $(\sigma, \boldsymbol{\mu})$ is consistent. Since $\mathrm{supp}(\sigma_i(\cdot | h)) \subseteq \arg\max_{a_i \in A_i(h)} \mathrm{E}[u_i | h, a_i]$ for all $i$ and $h \in H \backslash Z$, no player can profit from pure or randomized one-shot-deviations from $\sigma$. Since $\Gamma$ is finite, the one-shot-deviation principle applies, implying that $\sigma$ is subgame perfect in $\Gamma$. Now suppose that $(\sigma, \boldsymbol{\mu})$ is consistent; if $\sigma$ is a also a subgame perfect equilibrium of $\Gamma$ then the sequential rationality condition (3) of Definition 8 is satisfied; therefore $(\sigma, \boldsymbol{\mu})$ is a psychological sequential equilibrium.∎

**Corollary 11** *Suppose that $\Gamma$ is a standard game ($u_i$ depends only on $z$ for all $i$). Then, for any behavioral profile $\sigma$, $(\sigma, \boldsymbol{\mu})$ is a psychological sequential equilibrium for some $\boldsymbol{\mu}$ if and only if $\sigma$ is a subgame perfect (hence sequential) equilibrium of $\Gamma$.*

## 4.3  Examples

We illustrate this definition with three examples. The first example is a simultaneous move game which serves to illustrate how we can (in essence) reproduce the spirit of a leading example of GPS'. We then consider two psychological versions of the Trust Game, which connect back to some of the key notions previously highlighted in section 2.

    *Equilibrium beliefs in the Bravery Game.*
    The Bravery Game is a numerical example used in GPS (p. 66) to show that a psychological game may have multiple, isolated mixed strategy equilibria even if there is only one active player, which is impossible in standard games. We consider a modified version of their game to illustrate, in a very simple case, our definition of equilibrium in beliefs. The game is as follows.

There is only one active player: $A_1 = \{bold, timid\}$, $A_2 = \{Wait\}$ (2 is inactive). Since player 2 is inactive we can ignore his payoff function, but his beliefs do matter. Player 1 is concerned about what player 2 thinks about him. Acting boldly is dangerous, but it is worthwhile if player 2 expects player 1 to act boldly. GPS model the situation with a payoff function of the form $\overline{u}_1 : A \times \overline{\mathbf{M}}_1 \to \mathbb{R}$. Specifically, let $\alpha := \mu_2^1(bold|h^0)$ denote the first-order belief of player 2 about player 1 (a random variable from 1's point of view), and let $\beta := \mathrm{E}_{\cdot 1}[\alpha|h^0]$ denote (a feature of) the second-order beliefs of Player 1. The payoff function considered by GPS is

$$\overline{u}_1(a_1, \overline{\boldsymbol{\mu}}_1) = \left\{ \begin{array}{ll} 2 - \beta, & \text{if } a_1 = bold \\ 3(1 - \beta), & \text{if } a_1 = timid \end{array} \right.$$

We modify the payoff function of GPS so that it has the form $u_1 : A \times \mathbf{M} \to \mathbb{R}$. Specifically, we let

$$u_1(a_1, \boldsymbol{\mu}) = \left\{ \begin{array}{ll} 2 - \alpha, & \text{if } a_1 = bold \\ 3(1 - \alpha), & \text{if } a_1 = timid \end{array} \right.$$

Clearly, the expectation of $u_1$ given $a_1$ and player 1's second order belief $\beta$ is $\overline{u}_1$. It is easily checked that there are three equilibria: $\beta = \alpha = 1$, $\beta = \alpha = 0$ and $\beta = \alpha = \frac{1}{2}$.[29]

*Trust Game with Guilt Aversion*

Consider the psychological game $\Gamma_4$ (or equivalently game $\Gamma_2$). Recall that $\alpha = \alpha(\mu_1^1)$ is the probability that Ann assigns to strategy *'Share if Trust'* at the beginning of the game, and $\beta = \int \alpha(\mu_1^1)\mu_2^2(d\mu_1^1|Trust)$ is the relevant summary statistic of the second-order beliefs of Bob. We also let $\tau = \mu_2^1(Trust|h^0)$ denote the initial first-order belief of Bob. In this game an assessment is summarized by $(\tau, \alpha, \beta)$, where $(\tau, \alpha)$ corresponds to a behavior strategy profile. The indifference condition for Bob is $\beta = \frac{2}{5}$, the indifference condition for Ann is $\alpha = \frac{1}{2}$; consistency yields $\alpha = \beta$.

The game has three equilibrium assessments: $\tau = \alpha = \beta = 1$ (trust), $\tau = \alpha = \beta = 0$ (no trust), and $\tau = 0$, $\alpha = \beta = \frac{2}{5}$ (insufficient trust). Note that only the first equilibrium is consistent with forward induction reasoning (as described in section 2, and further elaborated on in subsection 5.1 below).

---

[29]These are essentially the same equilibria obtained by GPS. But they allow for explicit randomization; thus the first-order beliefs of Player 2 are degenerate on the equilibrium (mixed) strategy of Player 1, and higher-order beliefs of each player are degenerate on the equilibrium lower-order beliefs of the other player.

*Trust Game with Reciprocity*

Rabin (1993) illustrates how modeling reciprocity may require belief-dependent utilities, as kindness and perceived kindness depend on beliefs. Example 2 in the introduction provides an illustration. Our framework is adequate for modeling reciprocity in extensive games. To support this claim, we show how the essence of one particular form of reciprocity theory for extensive games, Dufwenberg & Kirchsteiger's, can be captured in an example which builds on $\Gamma_1$: Let $\alpha, \beta$,and $\tau$ be defined as in the previous example. The key tenets of the theory concern player $i$'s kindness to player $j$ ($K_{ij}$), and $i$'s belief in $j$'s kindness to $i$ ($\hat{K}_{iji}$). At each history, player $i$ maximizes utility defined by the sum of material payoffs (as in $\Gamma_1$) and reciprocity payoffs equal to $\theta_i \times K_{ij} \times \hat{K}_{iji}$, where $\theta_i$ is a constant measuring $i$'s sensitivity to reciprocity. Assume that Ann's sensitivity is $\theta_1 = 0$ and that Bob's sensitivity is $\theta_2 = \frac{4}{3}$. One can show that all relevant kindness notions can be reproduced in our framework and notation; in particular we need the following for Bob:

- Bob's kindness following *Trust* ($= K_{21}$) $= -1$ or $1$, for choices *Grab* and *Share*, respectively,

- Bob's belief in Ann's kindness following *Trust* ($= \hat{K}_{212}$) $= \frac{3}{2} - \beta$.

In $\Gamma_6$, we depict the relevant utilities as conceived by the players when they move (since Bob is not active at the root we put no utility for him following *Don't*):[30]
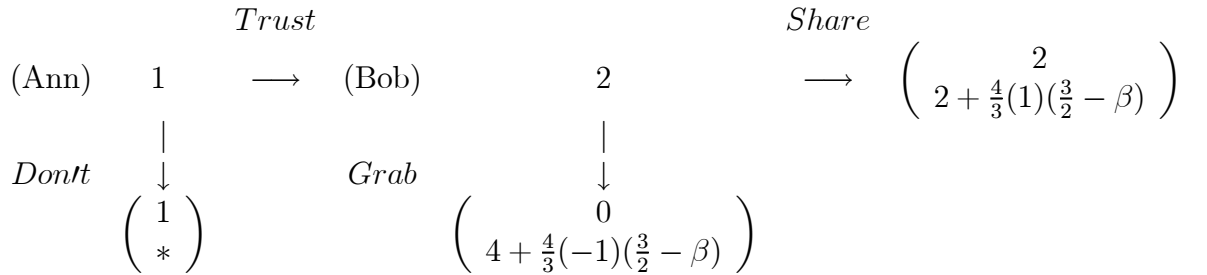


**Figure 6.** Trust Game $\Gamma_6$ with Reciprocity Payoffs.

Applying Definition 8, one sees that $\Gamma_6$ has a unique equilibrium assessments: $\tau = 1$, $\alpha = \beta = \frac{3}{4}$. No 'pure' PSE exists.[31] The prediction matches Dufwenberg & Kirchsteiger's.

---

[30]As in the Trust Game with guilt aversion, we can replace Bob's conditional second-order belief $\beta$ with Ann's initial first-order belief $\alpha$ in Bob's payoffs, and so obtain a strategically equivalent game.

[31]In any PSE we have $\alpha = \beta$. If $\theta_2 = \frac{4}{3}$, the indifference condition for Hugh yields $\beta = \frac{3}{4}$. If $\alpha = \beta < \frac{3}{4}$ then $\hat{K}_{212}$ shoots up, so Hugh prefers *Share* to *Grab*, which in PSE would imply $\alpha = \beta = 1$, ... a contradiction. If $\alpha = \beta > \frac{5}{6}$ then $\lambda_{212}$ goes down, so Hugh prefers *Grab* to *Share*, implying $\alpha = \beta = 0$, ... another contradiction.

# 5 Interactive Epistemology

We argued in section 2 that alternatives to equilibrium analysis are even more worth exploring for psychological games than for standard games. Fortunately, the definition of a psychological game provides us with all the ingredients to analyze strategic reasoning by means of interactive epistemology assumptions, that is, assumptions about players' rationality and what they believe about each other at any point of the game. In this section, we show how to express these assumptions in the language of events and belief operators (5.1), analyze a notion of rationalizability for dynamic psychological games under the simplifying assumption of 'own-strategy independence' (5.2), and discuss this assumption and how to deal with its removal (5.3).

## 5.1 States of the world, events, and belief operators

A state of the world is a complete specification of what the players would do and believe at each history of the game. Note the subjunctive conditional: game-theoretic analysis does not only concern the actual path of actions and beliefs, it must also consider how players would react (in terms of beliefs and choices) to histories that do not actually occur at the true state. The state of a player is therefore given by his strategy and his hierarchy of cps', $(s_i, \boldsymbol{\mu}_i)$. The set of states for player $i$ is denoted by $\Omega_i = S_i \times \mathbf{M}_i$, and the set of *states of the world* is $\Omega = \prod_{i=1}^{n} \Omega_i$. We let $\Omega_{-i} = \prod_{j \neq i} \Omega_j$ denote the set of possible states of $i$'s opponents. With a slight abuse of notation we often write $\Omega = \Omega_i \times \Omega_{-i}$ with typical element $\omega = (\omega_i, \omega_{-i})$.

An *event* is a (Borel) subset $E \subseteq \Omega$; its complement is denoted $\neg E = \Omega \backslash E$. An event about player $i$ is any set of states $E = E_i \times \Omega_{-i}$, where $E_i$ is a Borel subset of $\Omega_i$. We let $\mathcal{E}_i$ denote the family of events about $i$. Events about the opponents of player $i$ are similarly defined, and the collection of such events is denoted by $\mathcal{E}_{-i}$.

We often use brackets to denote specific events. In particular, for any function $\mathbf{x} : \Omega \to X$ and value $x^* \in X$, we use the notation $[\mathbf{x} = x^*] := \{\omega : \mathbf{x}(\omega) = x^*\}$. When $\mathbf{x}$ is understood, we simply write $[x^*]$. For example, $[s_i^*] = \{(s_i, \boldsymbol{\mu}_i, \omega_{-i}) : s_i = s_i^*\} \in \mathcal{E}_i$ is the event "$i$ plays $s_i^*$", where it is understood that $\mathbf{x}$ is the projection function on $S_i$, that is $\mathbf{x}(s_i, \boldsymbol{\mu}_i, \omega_{-i}) = s_i$. Similarly, $[h] = \prod_{i \in N} S_i(h) \times \mathbf{M}_i$ is the event that history $h$ occurs.

We follow both GPS and Battigalli & Siniscalchi (2002) in disregarding players' beliefs about themselves. At state $\omega = (s_i, \boldsymbol{\mu}_i, \omega_{-i})$, player $i$ would believe event $E = \Omega_i \times E_{-i} \in \mathcal{E}_{-i}$ conditional on history $h$ with probability $f_{i,h}(\boldsymbol{\mu}_i)(E_{-i})$ (cf. subsection 3.2). Thus $\{(s_i, \boldsymbol{\mu}_i, \omega_{-i}) : f_{i,h}(\boldsymbol{\mu}_i)(E_{-i}) = 1\}$ is the event "player $i$ would believe $E$ conditional on $h$". $E$ itself may be an

event concerning the beliefs of $i$'s opponents.

We use belief operators to represent events about interactive beliefs in a terse form: a *belief operator* for player $i$ is a mapping with domain $\mathcal{E}_{-i}$ and range $\mathcal{E}_i$. For any given history $h \in H$, the *h-conditional belief operator* for player $i$ is defined as follows:

$$\forall E = \Omega_i \times E_{-i} \in \mathcal{E}_{-i}, \, \mathrm{B}_{i,h}(E) = \{(s_i, \boldsymbol{\mu}_i, \omega_{-i}) : f_{i,h}(\boldsymbol{\mu}_i)(E_{-i}) = 1\}.$$

Note that $h$ may be counterfactual at $\omega$, because the strategies played at $\omega$ may not induce history $h$; in this case "$i$ would believe $E$ conditional on $h$" is a counterfactual statement about $i$'s beliefs (also called 'epistemic counterfactual'). Clearly, $\mathrm{B}_{i,h}(E) \in \mathcal{E}_i$.[32] Note that each $\mathrm{B}_{i,h}(\cdot)$ satisfies monotonicity [$E \subseteq F$ implies $\mathrm{B}_{i,h}(E) \subseteq \mathrm{B}_{i,h}(F)$] and conjunction [$\mathrm{B}_{i,h}(E \cap F) = \mathrm{B}_{i,h}(E) \cap \mathrm{B}_{i,h}(F)$]. Furthermore $\mathrm{B}_{i,h}(E) = \mathrm{B}_{i,h}(E \cap [h])$ because $i$ always believes what he observes.

The basic event we are interested in is players' rationality. To simplify our definition of rationality, for the time being we assume *own-strategy independence* (we discuss this assumption further in subsection 5.3), meaning that psychological preferences can be represented by a utility function of the form

$$u_i : Z \times \mathbf{M} \times S_{-i} \to \mathbb{R}. \tag{4}$$

The expectation of $u_i$ conditional on $h$, given $\boldsymbol{\mu}_i$ and $s_i$ is

$$\mathrm{E}_{s_i, \,_i}[u_i|h] = \int_{S_{-i} \times \mathbf{M}_{-i}} u_i(\zeta(s_i, s_{-i}), \boldsymbol{\mu}_i, \boldsymbol{\mu}_{-i}, s_{-i}) f_{i,h}(\boldsymbol{\mu}_i)(ds_{-i}, d\boldsymbol{\mu}_{-i}).$$

Following Pearce (1984), Rubinstein (1991), Reny (1992), and others, we take the point of view that the basic notion of rationality in extensive form games refers to plans-of-action rather than strategies; a rational player does not have to plan in advance what he would do if he deviated from his own plan. We say that player $i$ is *rational* at state $(s_i, \boldsymbol{\mu}_i, \omega_{-i})$ iff $s_i$ maximizes $i$'s conditional expected payoff, given belief hierarchy $\boldsymbol{\mu}_i$, conditional on each history allowed by $s_i$; more formally, let $H_i(s_i^*) = \{h \in H \backslash Z : s_i^* \in S_i(h)\}$ denote the set of non-terminal histories allowed by a fixed strategy $s_i^*$, then we require $s_i \in r(\boldsymbol{\mu}_i)$ where

$$r_i(\boldsymbol{\mu}_i) = \left\{ s_i^* : \forall h \in H(s_i^*), \, s_i^* \in \arg \max_{s_i \in S_i(h)} \mathrm{E}_{s_i, \,_i}[u_i|h] \right\} \tag{5}$$

The event that player $i$ is rational is $R_i = \{(s_i, \boldsymbol{\mu}_i, \omega_{-i}) : s_i \in r_i(\boldsymbol{\mu}_i)\}$. The assumption of own-strategy independence guarantees that $r_i(\boldsymbol{\mu}_i)$ can

---

[32]For any Borel set $\Omega_i \times E_{-i}$, $\mathrm{B}_{i,h}(\Omega_i \times E_{-i})$ is also a Borel set because the $h$-coordinate belief function $f_{i,h}$ is continuous (see Lemma 3).

be obtained via a backward induction algorithm and $R_i$ is a well-defined nonempty event (see the proof of Lemma 15 in the appendix).

To illustrate how these concepts can be used, without assuming equilibrium, we re-examine two psychological versions of the Trust Game. As a matter of notation, we have to distinguish the event "Bob shares", which in this extensive form implies that "Ann trusts Bob," from the event "Bob would share if Ann trusted Bob" which is a subjunctive conditional, logically independent on whether Ann trusts Bob or not. Similar considerations hold for the other action, *Grab*. We use **bold** letters to denote the subjunctive conditionals (which in this particular case correspond to strategies of Bob), as in [**Share**] and [**Grab**].

Consider the Trust Game with guilt aversion $\Gamma_4$. The game can be solved by forward induction reasoning: it is rational for Ann to trust Bob only if she assigns at least 50% probability to strategy **Share**, *i.e.* only if $\alpha \geq \frac{1}{2}$, where $\alpha : \mathbf{M}_1 \rightarrow \mathbb{R}$ is the random variable defined by $\alpha(\boldsymbol{\mu}_1) = \mu_1^1(\mathbf{Share}|h^0)$.[33] If Bob believes in Ann's rationality when he has to move (even if he is 'surprised'), he infers from Ann's action *Trust* that $\alpha \geq \frac{1}{2}$. Therefore $\beta \geq \frac{1}{2}$, where $\beta : \mathbf{M}_2 \rightarrow \mathbb{R}$ is the random variable defined by $\beta(\boldsymbol{\mu}_2) = \int \alpha(\boldsymbol{\mu}_1) f_{2,Trust}(\boldsymbol{\mu}_2)(d\boldsymbol{\mu}_1)$. His rational response is to share. If Ann anticipates Bob's reasoning she trusts him.

The formal counterpart of this argument is as follows (the events listed are nonempty; we rely on the monotonicity of the belief operators):

$$R_1 = \left\{ (s_1, \boldsymbol{\mu}_1, \omega_2) : \alpha(\boldsymbol{\mu}_1) > \frac{1}{2} \Rightarrow s_1 = Trust, \alpha(\boldsymbol{\mu}_1) < \frac{1}{2} \Rightarrow s_1 = Don't \right\}$$

$$R_2 = \left\{ (\omega_1, s_2, \boldsymbol{\mu}_2) : \beta(\boldsymbol{\mu}_2) > \frac{2}{5} \Rightarrow s_2 = \mathbf{Share}, \beta(\boldsymbol{\mu}_2) < \frac{2}{5} \Rightarrow s_2 = \mathbf{Grab} \right\},$$

$$R_1 \cap [Trust] \subseteq \left[ \alpha \geq \frac{1}{2} \right],$$

$$\mathrm{B}_{2,Trust}(R_1) = \mathrm{B}_{2,Trust}(R_1 \cap [Trust]) \subseteq \mathrm{B}_{2,Trust}\left( \left[ \alpha \geq \frac{1}{2} \right] \right) \subseteq \left[ \beta \geq \frac{1}{2} \right],$$

$$R_2 \cap \mathrm{B}_{2,Trust}(R_1) \subseteq R_2 \cap \left[ \beta \geq \frac{1}{2} \right] \subseteq [\mathbf{Share}],$$

$$R_1 \cap \mathrm{B}_{1,h^0}(R_2 \cap \mathrm{B}_{2,Trust}(R_1)) \subseteq R_1 \cap [\alpha = 1] \subseteq [Trust].$$

Now consider the Trust Game with Reciprocity $\Gamma_6$ in Figure 6 (or the equivalent version with $\beta$ replaced by $\alpha$). If one tries to analyze that game

---

[33]In some formulas, we have to make explicit the dependence of random variable $\alpha$ on the state of the world. The same holds for random variable $\beta$.

without an equilibrium supposition, one is at loss for predictive power: with the given sensitivity parameter $\theta_2 = \frac{4}{3}$, Bob's best response depends on whether $\beta$ is below or above the threshold $\left(\frac{3}{2} - \frac{1}{\theta_2}\right) = \frac{3}{4}$. This cannot be resolved by forward induction reasoning, which yields (as explained in sub-section 4.3) $\beta \geq \frac{1}{2}$.

However, if one modifies the game using other values of $\theta_2$ one can draw clear conclusions merely on the basis of backward induction: if $\theta_2 < \frac{2}{3}$, Bob's best response is *Grab* independently of $\beta$, thus $R_2 \subseteq [\textbf{Grab}]$ and $R_1 \cap B_{1,h^0}(R_2) \subseteq [Don't]$; on the other hand, if $\theta_2 > 2$, $R_2 \subseteq [\textbf{Share}]$ and $R_1 \cap B_{1,h^0}(R_2) \subseteq [Trust]$. Furthermore, a subtle issue arises when $\frac{2}{3} < \theta_2 < 1$. In this case backward induction cannot pin down Bob's best response, which is *Grab* if $\beta \geq \left(\frac{3}{2} - \frac{1}{\theta_2}\right)$; but a forward induction yields $\beta \geq \frac{1}{2}$. This puts an *upper bound* on how kind Bob believes that Ann is,[34] and with $\frac{2}{3} < \theta_2 < 1$ the best response is *Grab*. Formally, with these parameter values $R_2 \cap B_{1,Trust}(R_1) \subseteq [\textbf{Grab}]$, $B_{1,h^0}(R_2 \cap B_{1,Trust}(R_1)) \subseteq [\alpha = 0]$ and $R_1 \cap B_{1,h^0}(R_2 \cap B_{1,Trust}(R_1)) \subseteq [Don't]$.

One can show that, by contrast, the PSE prediction entails that $0 < \alpha = \beta = \left(\frac{3}{2} - \frac{1}{\theta_2}\right) < \frac{1}{2}$, $\tau = 0$. Thus, PSE and forward induction reasoning yield the same path, but very different predictions about how Bob would revise his beliefs off that path.

## 5.2 Rationalizability

The analysis of $\Gamma_4$ and $\Gamma_6$ in section 5.1 shows that PSE (like sequential equilibrium in standard games) need not be consistent with forward-induction reasoning. Here we provide the tools to perform a forward-induction analysis of general psychological games. Following Battigalli & Siniscalchi (2002), we first define a '*strong belief* operator' $SB_i$ as follows: $SB_i(\emptyset) = \emptyset$ and

$$\forall E \in \mathcal{E}_{-i} \backslash \{\emptyset\}, \ SB_i(E) = \bigcap_{[h] \cap E \neq \emptyset} B_{i,h}(E).$$

In words, $SB_i(E)$ is the event "player $i$ would believe $E$ conditional on every history that does not contradict $E$".[35] For example, $SB_i([s_j])$ is the event "*player $i$ would believe that player $j$ plays strategy $s_j$ at each history $h$ allowed by $s_j$*".

---

[34] The higher $\beta$, the more Bob believes that Ann's choice to trust him is self-interested.

[35] $SB_i(\cdot)$ is not a monotone operator, and it satisfies only a weak form of conjunction $[SB_i(E) \cap SB_i(F) \subset SB_i(E \cap F)]$. For more on this, see Battigalli & Siniscalchi (2002).

We will be interested in events of the form $\text{SB}_i(R_{-i} \cap E)$, where $R_{-i} = \bigcap_{j \neq i} R_j$ is the event that $i$'s opponents are rational and $E$ is either $\Omega$ or some event concerning beliefs, and we will consider assumptions like "everybody strongly believes that the opponents are rational." To write this in a simple form, we define a *mutual strong belief* operator. Let $\mathcal{E}$ denote the collection of events of the form $E = \bigcap_{i \in N} E_i$ ($E_i \in \mathcal{E}_i$). For example, $R = \bigcap_{i \in N} R_i \in \mathcal{E}$. For each $E \in \mathcal{E}$, the event "there is mutual strong belief in $E$" is defined by

$$\text{SB}(E) = \bigcap_{i \in N} \text{SB}_i \left( \bigcap_{j \neq i} E_j \right). \text{ Note that } \text{SB}(E) \in \mathcal{E}.$$

We explore the consequences of the following assumptions:

(0) each player is rational $[=R]$,

(1) mutual strong belief in (0) $[=\text{SB}(R)]$,

(2) mutual strong belief in (0) & (1) $[=\text{SB}(R \cap \text{SB}(R))]$,

(3) mutual strong belief in (0), (1) & (2) $[=\text{SB}(R \cap \text{SB}(R \cap \text{SB}(R)))]$,

and so on.... Such assumptions are more easily expressed with formulas if we introduce an auxiliary 'correct strong belief' operator:

$$\forall E \in \mathcal{E}, \ \text{CSB}(E) = E \cap \text{SB}(E)$$

The conjunction of assumptions $(0)$-$(k)$ above corresponds to the event $\text{CSB}^k(R)$, where for any $E \in \mathcal{E}$, $\text{CSB}^0(E) = E$ and $\text{CSB}^k(E) = \text{CSB}(\text{CSB}^{k-1}(E))$.[36] Rationalizability is defined by considering the limiting intersection for all $k$:

**Definition 12** *A state of the world $\omega$ is* rationalizable *if $\omega \in \bigcap_{k \geq 0} \text{CSB}^k(R)$. A strategy is rationalizable if it is part of a rationalizable state of the world.*

Battigalli & Siniscalchi (2002) show that the strategies consistent with event $\text{CSB}^k(R)$ in standard games are those that survive the first $k+1$ steps of Pearce's (1984) extensive-form rationalizability procedure. This explains the terminology of Definition 12.[37] To illustrate the concept, we can note that it captures the forward induction solution of game the Trust Game with guilt aversion (either $\Gamma_2$ or $\Gamma_4$). However, that conclusion requires only two layers of mutual correct strong belief, since the psychological forward induction solution in $\Gamma_2$ or $\Gamma_4$ obtains at all states $\omega \in \text{CSB}^2(R)$.

To illustrate the full power of Definition 12, we therefore analyze a *Generalized Trust Game with guilt aversion*, reminiscent of the Ben Porath &

---

[36]For example, (0) & (1) is $R \cap \text{SB}(R) = \text{CSB}(R)$, (0) & (1) & (2) is $R \cap \text{SB}(R) \cap \text{SB}(R \cap \text{SB}(R)) = \text{CSB}^2(R)$, etc.

[37]Alternative notions of rationalizability for extensive-form games have been explored. See, *e.g.*, the references in Battigalli & Siniscalchi (1999, 2002).

Dekel (1992) money-burning game: Ann can either (evenly) distribute the total surplus of \$2 or reinvest it in one out of $L$ projects. Project $\ell = 1, ..., L$ yields \$$2\left(1 + \frac{\ell}{L}\right)$. Bob controls the distribution of this larger surplus and can either *Grab* or (evenly) *Share*. We let $Trust_\ell$ denote the action of investing in project $\ell$, and $\mathbf{Share}_\ell$ denote the conditional choice of sharing if Ann invests in project $\ell$. Let $\alpha_\ell(\boldsymbol{\mu}_1) = \mu_1^1(\mathbf{Share}_\ell | h^0)$ and $\beta_\ell(\boldsymbol{\mu}_2) = \int \alpha_\ell(\mu_1^1)\mu_2^2(d\alpha_\ell(\mu_1^1)|Trust_\ell)$. As before we assume that Ann's utility is her material payoff, whereas Bob is averse to guilt. Applying the guilt formula of subsection 3.3, the players' utilities if Ann invests in project $\ell$ are given by

$$
\begin{aligned}
u_i(Trust_\ell, Share) &= \left(1 + \frac{\ell}{L}\right), \ i = 1, 2, \\
u_1(Trust_\ell, Grab) &= 0, \\
u_2(Trust_\ell, Grab) &= 2\left(1 + \frac{\ell}{L}\right) - \theta_2 \alpha_\ell \left(1 + \frac{\ell}{L}\right),
\end{aligned}
$$

where as before $\theta_2$ is Bob's sensitivity to guilt. Bob (strictly) prefers to share if and only if $\theta_2 \beta_\ell > 1$.

Note that if $L = 1$ we obtain the material payoff game $\Gamma_1$, and setting $\theta = \frac{5}{2}$ we obtain the psychological game $\Gamma_4$. Note also that the forward induction argument used to solve $\Gamma_4$ (captured by 2 iterations of the CSB operator) would work for every $\theta_2 > 2$, but would not work for lower values of $\theta_2$. On the other hand, in the Generalized Trust Game with guilt aversion rationalizability yields the efficient sharing outcome also for much lower values of $\theta_2$:

**Proposition 13** *In the Generalized Trust Game with guilt aversion, if $\theta_2 > 1 + \frac{1}{L}$ then, for every rationalizable state $(s_1, \boldsymbol{\mu}_1, s_2, \boldsymbol{\mu}_2)$, $s_1 = Trust_L$, $s_2 = (\mathbf{Share}_\ell)_{\ell=1}^L$, $\alpha_\ell(\boldsymbol{\mu}_1) = \beta_\ell(\boldsymbol{\mu}_2) = 1$ ($\ell = 1, ..., L$).*

**Proof.** If Ann chooses project $\ell$ she signals that $\alpha_\ell \geq \frac{L}{L+\ell}$, because she can obtain \$1 by not investing. By forward induction [event $\mathrm{CSB}(R)$], $\beta_\ell \geq \frac{L}{L+\ell}$ and Bob shares if $\theta_2 \frac{L}{L+\ell} > 1$, or equivalently $\theta_2 > 1 + \frac{\ell}{L}$. Since $\theta_2 > 1 + \frac{1}{L}$, event $\mathrm{CSB}^2(R)$ implies that Ann can (and will) secure the payoff $1 + \frac{\hat{\ell}(2)}{L} > 1$, where $\hat{\ell}(2)$ is largest $\ell$ such that $\theta_2 > 1 + \frac{\ell}{L}$. Now suppose that $\mathrm{CSB}^{2k}(R)$ imply that Ann can (and will) secure the payoff $1 + \frac{\hat{\ell}(2k)}{L}$ because she correctly believes that Bob would share the yield of each project $\ell = 1, ..., \hat{\ell}(2k)$. If $\hat{\ell}(2k) = L$, we are done. Otherwise, suppose that Ann chooses project $\ell > \hat{\ell}(2k)$. This signals that $\alpha_\ell \geq \frac{L+\hat{\ell}(2k)}{L+\ell}$. Event

$\text{CSB}^{2k+1}(R)$ implies that $\beta_\ell \geq \frac{L+\hat{\ell}(2k)}{L+\ell}$ and Bob shares if $\theta_2 \frac{L+\hat{\ell}(2k)}{L+\ell} > 1$, or equivalently $\theta_2 > 1 + \frac{\ell - \hat{\ell}(2k)}{L+\hat{\ell}(2k)}$. Hence $\text{CSB}^{2k+2}(R) = \text{CSB}^{2(k+1)}(R)$ implies that Ann can (and will) secure the payoff $1 + \frac{\hat{\ell}(2(k+1))}{L}$, where $\hat{\ell}(2(k+1)) = \max\left\{\ell \in \{1, ..., L\} : \theta_2 \frac{L+\hat{\ell}(2k)}{L+\ell} > 1\right\}$. Since $1 + \frac{1}{L} \geq 1 + \frac{1}{L+\hat{\ell}(2k)}$, assumption $\theta_2 > 1 + \frac{1}{L}$ implies that function $\hat{\ell}(2k)$ is well-defined and strictly increasing until it attains its maximum, $L$. The thesis easily follows. ∎

The following theorem shows that our extension of Pearce's solution concept to psychological games is well behaved.

**Theorem 14** *If psychological payoffs are continuous functions and satisfy own-strategy independence (4), the set $\bigcap_{k \geq 0} \text{CSB}^k(R)$ of rationalizable states is nonempty and compact.*

**Proof.** By definition

$$\text{CSB}^{k+1}(R) = \text{CSB}(\text{CSB}^k(R)) = \text{CSB}^k(R) \cap \text{SB}(\text{CSB}^k(R)) \subseteq \text{CSB}^k(R).$$

[**NOTE**]We prove below by induction that each element $\text{CSB}^k(R) = \bigcap_{\ell=0}^{k} \text{CSB}^\ell(R)$ of the nested sequence $\left\{\text{CSB}^k(R)\right\}_{k \geq 0}$ is closed and nonempty. Lemma 3 implies that $\Omega$ is compact; thus, the closed subset $\bigcap_{k \geq 0} \text{CSB}^k(R)$ is compact. Furthermore, the finite intersection property of compact spaces implies that $\bigcap_{k \geq 0} \text{CSB}^k(R) \neq \emptyset$.

The inductive argument relies on the following three preliminary results, which are proved in the appendix.

**Lemma 15** *If the payoff function of player $i$ has the form (4), correspondence $r_i : \mathbf{M}_i \twoheadrightarrow S_i$ is nonempty valued. If $u_i$ is also continuous, $r_i$ has a closed graph and $R_i$ is a nonempty closed set.*

**Lemma 16** *For every closed event $E \in \mathcal{E}$, $SB(E)$ is closed.*

**Lemma 17** *Let $\left\{E^\ell\right\}_{\ell=0}^{\ell=k}$ be a decreasing sequence of nonempty events in $\mathcal{E}$ ($\emptyset \neq E^k \subseteq E^{k-1} \subseteq ... \subseteq E^0$), then $\bigcap_{\ell=0}^{\ell=k} SB(E^\ell)$ is also nonempty.*

For notational convenience let $\text{CSB}^{-1}(E) = \Omega$. We prove by induction that, for each $k \geq 0$, $\text{CSB}^k(R)$ is nonempty closed and can be expressed as

$$\text{CSB}^k(R) = R \cap \left(\bigcap_{\ell=-1}^{k-1} \text{SB}\left(\text{CSB}^\ell(R)\right)\right).$$

*Basis step.* The statement is true for $k = 0$ because by Lemma 15 $\mathrm{CSB}^0(R) = R$ is nonempty closed, and by definition $R$ can be expressed as

$$R = R \cap \Omega = R \cap \mathrm{CSB}^{-1}(R)$$

*Inductive step.* Suppose that the statement is true for each $\ell = 0, ..., k$, then

$$
\begin{aligned}
\mathrm{CSB}^{k+1}(R) &= \mathrm{CSB}(\mathrm{CSB}^k(R)) = \mathrm{CSB}^k(R) \cap \mathrm{SB}(\mathrm{CSB}^k(R)) \\
&= R \cap \left( \bigcap_{\ell=-1}^{k-1} \mathrm{SB}\left(\mathrm{CSB}^\ell(R)\right) \right) \cap \mathrm{SB}(\mathrm{CSB}^k(R)) \\
&= R \cap \left( \bigcap_{\ell=-1}^{k} \mathrm{SB}\left(\mathrm{CSB}^\ell(R)\right) \right).
\end{aligned}
$$

By the inductive hypothesis each $\mathrm{CSB}^\ell(R)$ is nonempty and closed ($\ell = 0, ..., k$). By Lemma 16 also $\mathrm{SB}\left(\mathrm{CSB}^\ell(R)\right)$ is closed ($\ell = 0, ..., k$). $R$ is also closed (Lemma 15). Hence $\mathrm{CSB}^{k+1}(R)$ is closed. $\left\{\mathrm{CSB}^\ell(R)\right\}_{\ell=0}^{\ell=k}$ is a decreasing sequence of nonempty events in $\mathcal{E}$. Therefore Lemma 17 implies that $\bigcap_{\ell=-1}^{k} \mathrm{SB}\left(\mathrm{CSB}^\ell(R)\right) \neq \emptyset$. Pick any state $\omega = (s_i, \boldsymbol{\mu}_i)_{i \in N} \in \bigcap_{\ell=-1}^{k} \mathrm{SB}\left(\mathrm{CSB}^\ell(R)\right)$. Since the latter is just an event about beliefs, modifying the strategies in $\omega$ we obtain another state in the same event. By definition of $R$, $\prod_{i \in N} r_i(\boldsymbol{\mu}_i) \times \{\boldsymbol{\mu}_i\} \subseteq R$. By Lemma 15, $r_i(\boldsymbol{\mu}_i) \neq \emptyset$. We conclude that

$$\emptyset \neq \prod_{i \in N} r_i(\boldsymbol{\mu}_i) \times \{\boldsymbol{\mu}_i\} \subseteq R \cap \left( \bigcap_{\ell=-1}^{k} \mathrm{SB}\left(\mathrm{CSB}^\ell(R)\right) \right).$$

Hence $\mathrm{CSB}^{k+1}(R) \neq \emptyset$. This completes the proof of the inductive step and the proof of the theorem.■

## 5.3 Own-strategy dependence and dynamic consistency

So far in section 5 we have assumed own-strategy independence. This assumption allows us to apply standard dynamic programming techniques, and enables us to prove Theorem 13. The purpose of this subsection is to exhibit the array of problems one may run into if one does not insist on own-strategy independence, and to indicate what might be done to fix those problems.

We need to introduce some additional notation. For any $s_i^*$, $h \in H \backslash Z$ and $a_i \in A_i(h)$, let $S_i(h, s_i^*)$ denote the set of strategies that reach $h$ and coincide with $s_i^*$ at all histories that do not weakly follow $h$, let $a_i h s_i^*$ denote

the strategy that selects $a_i$ at $h$ and coincide with $s_i^*$ at all histories but $h$, and finally let $re(s_i^*)$ denote the set of strategies realization-equivalent to $s_i^*$.[38] Recall that $H(s_i^*)$ is the set non-terminal histories allowed by $s_i^*$.

Own-strategy *in*dependence implies the following equivalences: for all $\boldsymbol{\mu}_i$ and $s_i^*$

$$\forall h \ \in \ H(s_i^*), \ s_i^* \in \arg \max_{s_i \in S_i(h)} \mathrm{E}_{s_i, \ _i}[u_i|h] \tag{6}$$

$$\text{if and only if}$$

$$\forall h \ \in \ H(s_i^*), \ s_i^* \in \arg \max_{s_i \in S_i(h, s_i^*)} \mathrm{E}_{s_i, \ _i}[u_i|h] \tag{7}$$

$$\text{if and only if}$$

$$\exists \hat{s}_i \ \in \ re(s_i^*), \ \forall h \in H \backslash Z, \ \hat{s}_{i,h} \in \arg \max_{a_i \in A_i(h)} \mathrm{E}_{a_i h \hat{s}_i, \ _i}[u_i|h] \tag{8}$$

The equivalence between (6) and (7) holds because under own-strategy independence the behavior of $s_i$ at histories ruled out by $h$ cannot affect the expected payoff conditional on $h$. The equivalence between (7) and (8) is a version of the one-shot-deviation principle: it says that a strategy $s_i^*$ is immune to deviations in subtrees allowed by $s_i^*$ if and only if $s_i^*$ is realization-equivalent to some strategy $\hat{s}_i$ which is immune to one-shot deviations.
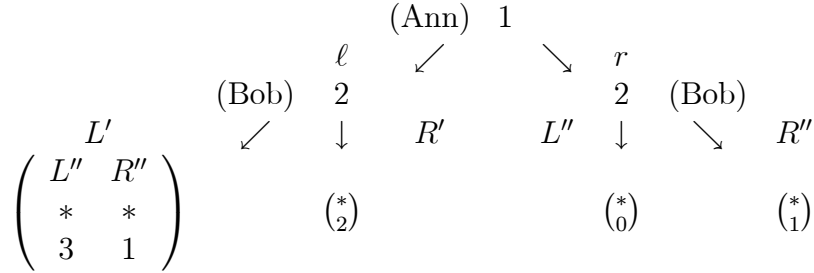


**Figure 7.** The game $\Gamma_7$

If we drop own-strategy independence these equivalences need not hold. Consider the game $\Gamma_7$. Here own-strategy independence fails because the payoff of Bob at terminal history $(\ell, L')$ depends on what he would have done had Ann chosen $r$ instead of $\ell$ (Ann's payoff is suppressed because it is not relevant to our argument). If bygones were not bygones, Bob would choose strategy $\mathbf{L'L''}$ at $\ell$ (we keep using use bold letters to denote conditional choices, to be distinguished from actual actions). But the best action after $r$ is $R''$. Therefore no strategy $s_2$ satisfies condition (6), whereas (7) and (8)

---

[38]$s_i$ is realization-equivalent to $s_i^*$ iff $\forall s_{-i}$, $\zeta(s_i, s_{-i}) = \zeta(s_i^*, s_{-i})$. Two realization-equivalent strategies allow the same set of nonterminal histories and select the same actions at such histories.

identify the strategy $\mathbf{R'R''}$ [(7) and (8) are equivalent in every game where player $i$ moves at most once in any play]. Therefore the equivalence between (6) and (7) [or (8)] fails.
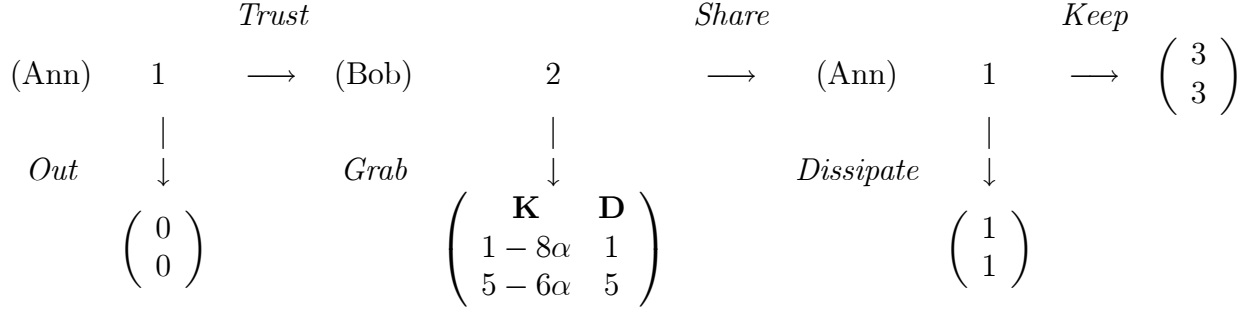


**Figure 8.** Modified Trust Game with Disappointment $\Gamma_8$.

Now consider game $\Gamma_8$. Here Ann's utility is affected by disappointment and Bob's utility by guilt.[39] As before $\alpha = \mu_1^1(\mathbf{Share}|h^0)$. Strategy (*Trust*, **Dissipate**) yields 1, and (*Trust*, **Keep**) yields $3\alpha + (1 - \alpha)(1 - 8\alpha)$. If $1 > 3\alpha + (1 - \alpha)(1 - 8\alpha)$, i.e. $0 < \alpha < \frac{3}{4}$, the equivalence between (7) and (8) fails: Indeed, no strategy satisfies (7) [or (6)], because (*Trust*, **Keep**) is the only maximizing strategy at history $h = (Trust, Keep)$, whereas (*Trust*, **Dissipate**) is the only maximizing strategy at the root. On the other hand, a strategy satisfying (8) can be found by backward induction: maximization at $h = (Trust, Share)$ yields **Keep**; given **Keep**, maximization at the root yields either *Out* (a reduced form strategy) or (*Trust*, **Keep**) according to whether $3\alpha + (1 - \alpha)(1 - 8\alpha) < 0$ or not. In the first case, that is, if $\frac{1}{4} < \alpha < \frac{1}{2}$, the strategy obtained by backward induction is *Out*. But Ann would be willing to pay up to \$1 to be able to commit to strategy (*Trust*, **Dissipate**). Thus, Ann preferences are dynamically inconsistent. Ann is unable to commit (otherwise this possibility should be modeled explicitly as a move in the extensive form). Therefore we argue that the relevant rationality condition is (8).

Note that if $\alpha = \frac{1}{3}$ and Bob has correct beliefs about $\alpha$ (at least in expectation), he is indifferent; hence $\alpha = \frac{1}{3}$ can be justified and *Out* is a sequential equilibrium outcome.[40]

It is natural to ask whether rationality condition (8) can always be satis-

---

[39]We use the formulas of subsection 3.3 with $\theta_1 = 4$ and $\theta_2 = 3$.

[40]*Out* is also a rationalizable outcome (given the appropriate notion of rationality), because *Trust* only signals that either $\alpha \le \frac{1}{4}$ or $\alpha \ge \frac{1}{2}$. Therefore rationalizability does not rule out either strategy of Bob, and $\frac{1}{4} \le \alpha \le \frac{1}{2}$ is a rationalizable first-order belief.

fied. The next example shows that the answer is No.

$$
\begin{array}{c}
\text{(Ann)} \quad 1 \\
\ell \quad \swarrow \qquad \searrow \quad r \\
\text{(Bob)} \quad 2 \qquad\qquad 2 \quad \text{(Bob)} \\
L' \quad \swarrow \quad \downarrow \quad R' \qquad L'' \quad \downarrow \quad \searrow \quad R'' \\
\begin{pmatrix} & L'' & R'' \\ & 2 & 2 \\ & 3 & 1 \end{pmatrix} \qquad \begin{pmatrix} 0 \\ 2 \end{pmatrix} \qquad\qquad \begin{pmatrix} & L' & R' \\ & 0 & 0 \\ & 0 & 3 \end{pmatrix} \qquad \begin{pmatrix} 2 \\ 1 \end{pmatrix}
\end{array}
$$

**Figure 9.** The Psychological Game $\Gamma_9$

Game $\Gamma_9$ is a modified version of game $\Gamma_9$. A strategy satisfies (8) if and only if it corresponds to a pure Nash equilibrium of the following game between two selves of Bob:

| Bob'\Bob" | $L''$ | $R''$ |
|-----------|-------|-------|
| $L'$      | 3,0   | 1,1   |
| $R'$      | 2,3   | 2,1   |

But this companion game has no pure equilibrium. On the other hand, solving for the indifference conditions we obtain a unique mixed equilibrium: $(\frac{2}{3}L' + \frac{1}{3}R', \frac{1}{2}L'' + \frac{1}{2}R'')$. Working backward, the best reply of Ann is $\ell$. This is the unique PSE of $\Gamma_9$ (recall that by Theorem 9 every continuous psychological game has a PSE).

This example shows that in order to make rationality possible in general psychological games we have to allow for the possibility that a player is uncertain about her own strategy, as we implicitly did in our analysis of sequential equilibrium in section 4. We can *explicitly* account for this uncertainty at the cost of some additional complexity. The first-order cps' of player $i$ must be defined as elements of $\Delta^H(S)$ rather than $\Delta^H(S_{-i})$. The construction of infinite hierarchies on beliefs goes through pretty much as in subsection 3.3. The additional difficulty is that now we have to deal explicitly with a player's beliefs about herself when we define rationality conditions. Following Battigalli & Siniscalchi (1999), one could include in the rationality condition that each player regards her strategy as stochastically independent of the strategies and beliefs of the opponents, but unlike that paper we have to allow for uncertainty about one's own strategy. Using a fixed-point argument one can show that, if $u_i$ is continuous, for every hierarchy of marginal cps' $\boldsymbol{\mu}_{i,-i}$ about the opponents of player $i$, there is a strategy $s_i$ and a cps about oneself $\mu^1_{i,i} \in \Delta^H(S_i)$ such that $(s_i, \mu^1_{ii}, \boldsymbol{\mu}_{i,-i})$ satisfies the one-shot deviation property. Using this property in the definition of rationality we can obtain a generalization of Theorem 14.

# 6 Discussion and Extensions

In this section we compare our framework with GPS (6.1) and provide three extensions, concerning incomplete information (6.2), imperfect observability of past actions (6.3), and a more general class of preferences whereby the move at different histories is controlled by different 'selves' of the same player with distinct 'local' utility functions (6.4).

## 6.1 Comparison with GPS

In section 2 we presented our framework as a generalization of GPS. But this is not literally true. The reason is twofold. First, GPS allow for imperfect information and chance moves. But, as we show below, these complications can be included in our framework. Second, GPS allow for explicit randomization whereas we exclude it. *Prima facie*, this difference may seem immaterial. Since GPS, like us, assume that players maximize their expected (psychological) utility given their beliefs, in their framework there is no incentive to randomize and it seems that the only role played by randomization is to guarantee the existence of equilibrium, a result that we obtain by looking at equilibrium in beliefs. But, unlike standard games, in psychological games there may be a difference (to a player's expected utility) between a belief that assigns probability one to the mixed strategy that, say, pick $a$ or $b$ with probability $\frac{1}{2}$, and the belief that assigns probability $\frac{1}{2}$ to $a$ and $\frac{1}{2}$ to $b$. These two beliefs are equivalent if psychological utility functions satisfy a linearity property. In most of the examples and applications of psychological games we are aware of this property is satisfied.

Let us now look at the version of GPS that *is* a special case of our framework: i.e. psychological games with utility functions of the form $u_i : Z \times \overline{\mathbf{M}}_i \to \mathbb{R}$, where $\overline{\mathbf{M}}_i$ is the space of infinite hierarchies of *initial* beliefs of player $i$, and first-order beliefs are probability measures over pure strategies of the opponents. How much is lost by restricting the analysis to such games? We have argued that many interesting phenomena such as sequential reciprocity, psychological forward induction and regret cannot be analyzed within this class of games. But we can prove a partial equivalence result. Suppose that the initial beliefs of others enter the utility function: $u_i : Z \times \prod_{j \in N} \overline{\mathbf{M}}_j \to \mathbb{R}$. Then there is a psychological game with utility functions $\overline{u}_i : Z \times \overline{\mathbf{M}}_i \to \mathbb{R}$ that has the same sequential equilibrium assessments as the former game.[41] This does not mean that in this class of games

---

[41]The intuition is relatively simple: each initial belief hierarchy $\overline{\boldsymbol{\mu}}_i$ induces a probability measure $\overline{f}_i(\overline{\boldsymbol{\mu}}_i) \in \Delta(S_{-i} \times \overline{\mathbf{M}}_{-i})$ which can be used to compute an expectation $\overline{u}_i(z, \overline{\boldsymbol{\mu}}_i)$

conditional-higher are immaterial. First, the equivalence result only concerns sequential equilibria, and we argued that the non-equilibrium analysis of psychological games is important. Second, our very definition of sequential equilibrium makes essential use of conditional beliefs.[42]

## 6.2 Incomplete information

Unless one models interaction within a family or amongst friends, it is probably not realistic to assume that players know one another's psychological propensities. Many of the examples we have looked at can be criticized on that ground. For example, in the analysis of game $\Gamma_2$ (or $\Gamma_4$) we assumed that Ann knows that Bob's sensitivity to guilt is $\theta_2 = \frac{5}{2}$, which may be a stretch.[43] Another reason to allow for incomplete information is that a player may care about the beliefs of others about some of his characteristics, which are not common knowledge, as in the models of Bernheim and Dufwenberg & Lundholm.

In order to extend the analysis of psychological games to include incomplete information, let $\theta = (\theta_0, \theta_1, ..., \theta_n)$ denote a vector of parameters that summarize all the payoff-relevant aspects of the game that are not common knowledge; $\theta_i$ is a component known to player $i$ only (such as his ability, or his sensitivity to certain psychological motivations), nobody knows $\theta_0$. It is common knowledge that $\theta$ belongs to a parameter space $\Theta = \Theta_0 \times \Theta_1 \times ... \times \Theta_n$. Elements of $\Theta$ are called *states of Nature*. We assume that $\Theta$ is a compact Polish space. We also assume for simplicity that players do not get more refined information about the state of Nature as the play unfolds, they only observe the actions chosen in previous stages of the game.

It is relatively easy to generalize our construction of the belief space in order to include beliefs about the state of Nature: replace $X^0_{-i} = S_{-i}$

---

of $u_i(z, \overline{\mu}_i, \cdot)$. Since in a consistent assessment there is 'common knowledge' of the hierarchical beliefs, no observation will make the players change their mind about the initial beliefs of the opponents, hence for any consistent assessment $u_i$ and $\overline{u}_i$ have the same set of maximizing actions at each history. (If there are simultaneous moves $u_i$ and $\overline{u}_i$ are fully equivalent, that is, they have the same best response correspondences.)

[42]If conditional beliefs were not in the language we would have to use an indirect approach similar to the one adopted by GPS: First define what a psychological Nash equilibrium is using the *ex ante* versions of Definitions 6 and 8. Then stipulate that a Nash equilibrium profile $\overline{\mu}$ is a sequential equilibrium if there is a behavioral strategy profile $\sigma$ that is a sequential equilibrium of the standard game with payoff functions $u_i(\cdot, \overline{\mu})$, and is such that each $\text{marg}_{S_i} \overline{\mu}_j^1 (j \neq i)$ is derived from $\sigma_i$ *via* Kuhn's transformation.

[43]There is ample evidence in psychology that emotional sensitivities differ among individuals. See Krohne (2003) for a general discussion. Tangney (1995) discusses guilt specifically.

with $X^0_{-i} = S_{-i} \times \Theta_{-i}$ in the construction of subsection 3.2, where $\Theta_{-i} = \Theta_0 \times ... \times \Theta_{i-1} \times \Theta_{i+1} \times ... \times \Theta_n$. Conditioning events for first-order beliefs now have the form $F = S_{-i}(h) \times \Theta_{-i}$ $(h \in H)$. Let $X^{k-1}_{-i}$ be the space of $(k-1)$-order uncertainty for player $i$; then we obtain the set of $k$-order cps'. Then we obtain the set of $k$-order cps' $\Delta^H(X^{k-1}_{-i})$, and the $k$-order uncertainty space $X^k_{-i} = X^{k-1}_{-i} \times \prod_{j \neq i} \Delta^H(X^{k-1}_{-j})$. Lemmata 2 and 3 are easily extended to this case. Therefore we obtain, for each $i \in N$, the space $\mathbf{M}_i$ of infinite hierarchies of cps' consistent with collective coherency, which is a compact Polish space homeomorphic to $\Delta^H(S_{-i} \times \Theta_{-i} \times \mathbf{M}_{-i})$.[44]

This is all we need to define the domain of psychological payoff functions, but it need not exhaust the description of the psychological game. Since states of Nature are exogenous, also players' hierarchies of *initial* beliefs about the state of Nature are exogenous,[45] and the model may specify assumptions about such exogenous beliefs. For example, one may assume that beliefs about $\theta$ are derived from a common prior $\rho \in \Delta(\Theta)$ and that this is common knowledge.[46] More sophisticated assumptions are allowed by Harsanyi's implicit representation of belief hierarchies by means of a $\Theta$-based type space (cf. Harsanyi, 1967-68, and Mertens & Zamir, 1985). Alternatively, assumptions about exogenous beliefs may be stated explicitly. Whatever these assumptions may be, they identify subspaces of hierarchies of cps' $\hat{\mathbf{M}}_i \subseteq \mathbf{M}_i$, $i = N$, which form the basis for the strategic analysis of the game. The analysis of rationalizability can be quite easily extended to this more general framework.[47] Psychological sequential equilibrium requires more care because the extension of the definition of consistency to general games of incomplete information is not obvious.

With this extended framework in place, we can regard (appropriately discretized versions of) the models of social conformity (Bernheim) and social respect (Dufwenberg & Lundholm) as psychological games with incomplete information. These models have a non-standard signaling game structure. There is only one active player, player 1, who has private information $\theta$ and chooses action $a$. Player 2 (who represents society) is passive, he only observes action $a$ and makes inferences about $\theta$. The payoff function of

---

[44]See Battigalli & Siniscalchi (1999).

[45]Posterior beliefs about the state of Nature are endogenous, because they are derived from joint beliefs about strategies (and beliefs of others) and state of Nature by conditioning on histories.

[46]Even in this simple case, we should distinguish the incomplete information situation from one where there is asymmetric information about chance moves, especially in the rationalizability analysis. On chance moves see the next subsection.

[47]A state of player $i$ is a triple $(s_i, \theta_i, \boldsymbol{\mu}_i)$; the definition of rationality is almost the same as in section 5, except that player $i$ takes into accout her knowledge $\theta_i$ of the state of Nature. See Battigalli & Siniscalchi (2002).

the active informed player has the form $u_1 : A \times \Theta \times \Delta^H(A \times \Theta) \to \mathbb{R}$ (where $H = \{h^0\} \cup A$). More specifically, there is a valuation function $v_1 : A \times \Theta \times \Delta(\Theta) \to \mathbb{R}$ such that

$$u_1(a, \theta, \mu_2) = v_1(a, \theta, \mu_2(\cdot|a)).$$

This means that Player 1 cares about the beliefs Player 2 will hold about private information $\theta$ conditional on his action $a$. (Note that this is an instance where *terminal* beliefs (of other players) affect payoffs.)

## 6.3  Imperfectly observable actions and chance moves

We chose to focus on games with observable actions and no chance moves for the sake of simplicity. But our concepts and results carry over to the more general case of games where past actions need not be perfectly observed and chance may play a role (as in GPS). We let $N = \{0, 1, ..., n\}$ where index 0 denotes the chance player, and let $\mathbf{H}_i$ be the partition of the the set of histories $H$ into information sets of player $i$ ($i \neq 0$).[48] Assume that perfect recall holds. Then it must be the case that the set of strategy profiles consistent with any information set $\mathbf{h}_i \in \mathbf{H}_i$ has the form $S(\mathbf{h}_i) = S_i(\mathbf{h}_i) \times S_{-i}(\mathbf{h}_i)$. We have to consider, for the first-order beliefs of player $i$, the collection of conditioning events $\{F_i : F_i = S_{-i}(\mathbf{h}_i), \mathbf{h}_i \in \mathbf{H}_i\}$. Let $X_{-i}^{k-1}$ be the space of $(k-1)$-order uncertainty for player $i$; then we obtain the set of $k$-order cps' $\Delta^{\mathbf{H}_i}(X_{-i}^{k-1})$, and the $k$-order uncertainty space $X_{-i}^k = X_{-i}^{k-1} \times \prod_{j \neq i, 0} \Delta^{\mathbf{H}_j}(X_{-j}^{k-1})$. The resulting set of infinite hierarchies of cps' $\mathbf{M}_i$ is homeomorphic to $\Delta^{\mathbf{H}_i}(S_{-i} \times \mathbf{M}_{-i})$. As in the case of incomplete information, the analysis of rationalizability is easily extended, while the definition of consistency in the sequential equilibrium analysis requires more care.[49]

## 6.4  Multi-self players and sequential reciprocity

We have seen in subsection 5.3 how own-strategy dependence may yield interesting instances of dynamic inconsistency. A more direct way to allow for

---

[48]Note that also terminal histories are partitioned into information sets because terminal beliefs are allowed to play a role.

[49]In particular, the definition of sequential equilibrium in beliefs requires an enlarged collection of conditioning events allowing a form of 'virtual conditioning' on the information sets of opponents; otherwise the behavior strategy of player $i$ cannot be derived from the conditional beliefs of player $j$. Furthermore, conditions (a) and (b) of Definition 6 must be extended to this more general framework. One way to do it, although not the most transparent, is to replace them with the topological condition originally used by Kreps & Wilson. Similar considerations apply to games with incomplete information.

dynamic inconsistencies is to adopt a multi-self approach and model a player preferences with an array of 'local' utility functions $(u_{i,h} : Z \times \mathbf{M} \times S \to \mathbb{R})_{h \in H \setminus Z}$. The sequential equilibrium analysis of section 4 applies to this extended framework almost *verbatim*.

We exemplify with reference to Dufwenberg & Kirchsteiger's reciprocity theory. We have already seen how our basic framework could reproduce their reciprocity theory in an example ($\Gamma_6$), but to handle general games one needs the multi-selves approach.[50] We consider two-person games for simplicity. Recall that *ceteris paribus* player $i$ likes to be kind toward $j$ if she believes that $j$ is kind toward her. Kindness depends on intentions. In particular, the *kindness of $j$ toward $i$*, $K_{ji}$, given $j$'s first-order belief $\nu \in \Delta(S_i)$ and $j$'s strategy $s_j$ is an increasing function of the difference between the expected material payoff of $i$ (given $s_j$ and $\nu$) and a belief-dependent "equitable payoff" $\pi_{ji}^e(\nu)$ that $j$ ascribes to $i$:

$$K_{ji}(s_j, \nu) = \sum_{s_i'} \nu(s_i') \pi_i(\zeta(s_i', s_j)) - \pi_{ji}^e(\nu).$$

The kindness of a player toward the co-player depends on his current first-order belief, which depends on the observed history. Therefore, for any fixed hierarchy of cps' $\boldsymbol{\mu}_j = (((\mu_j^1(\cdot|h))_{h \in H}, (\mu_j^2(\cdot|h))_{h \in H}, ...)$, the kindness of $j$ toward $i$ at history $h$ is $K_{ji}(s_j, \mu_j^1(\cdot|h))$, where $s_j \in S_i(h)$. One way to capture reciprocity motivations is to assume that at each history $h$ player $i$ maximizes the expected value of a linear combination of her material payoff and the product between her kindness at $h$ toward the opponent and the opponent's kindness at $h$ toward her, that is, $i$ at $h$ maximizes

$$\int_{S_j(h) \times \Delta(S_i(h))} \left[ \pi_i(\zeta(s_i, s_j)) + \theta_i K_{ij}(s_i, \mu_i^1(\cdot|h)) K_{ji}(s_j, \mu_j^1(\cdot|h)) \right] \mu_i^2(ds_j, d\mu_j^1(\cdot|h)|h).$$

This means that $i$'s preferences at history $h$ are represented by the psychological payoff function

$$u_{i,h}(z, \boldsymbol{\mu}, s) = \pi_i(z) + \theta_i K_{ij}(s_i, \mu_i^1(\cdot|h)) K_{ji}(s_j, \mu_j^1(\cdot|h)),$$

or equivalently by the function

$$\pi_i(z) + \theta_i K_{ij}(s_i, \mu_i^1(\cdot|h)) \hat{K}_{iji}(\mu_i^2(\cdot|h)),$$

where $\hat{K}_{iji}(\mu_i^2(\cdot|h)) = \int_{S_j(h) \times \Delta(S_i(h))} K_{ji}(s_j, \mu_j^1(\cdot|h)) \mu_i^2(ds_j, d\mu_j^1(\cdot|h)|h)$ is $i$'s belief in $j$'s kindness toward $i$. What we have here is, essentially, a reformulation of Dufwenberg & Kirchsteiger's model.[51]

---

[50]This is not to suggest that one could not conceive of a different sort of reciprocity theory, which would not require a multi-selves approach.

[51]There are inessential differences, concerning *e.g.* how beliefs and the domain of ma-

# 7    Concluding remarks

The utility of decision makers who are motivated by 'psychological' considerations such as reciprocity, guilt, social respect, or social conformity may depend directly on their beliefs (about others' choices, beliefs, or information). In a pioneering contribution, Geanakoplos, Pearce & Stacchetti point out that traditional game theory does not address such matters, and they present a model which does. However, their toolbox of 'psychological games' incorporates several restrictions that rule out many plausible forms of belief-dependent motivation. In particular, they cannot address the issue of how beliefs about the beliefs of others are revised as the play unfolds. We propose a more general framework, which allows updated higher-order beliefs, beliefs of others, planned strategies, and incomplete information to influence motivation. We develop new solution concepts, and provide examples and existence results.

The range of topics that to date have been explored in models of belief dependent motivation is limited. We propose that there are a variety of interesting forms of belief-dependent motivation waiting to be analytically explored. In his survey paper on "Emotions and Economic Theory", Elster (1998) argues that a key characteristic of emotions is that "they are triggered by beliefs" (p. 49). He discusses, inter alia, anger, hatred, guilt, shame, pride, admiration, regret, rejoicing, disappointment, elation, fear, hope, joy, grief, envy, malice, indignation, jealousy, surprise, boredom, sexual desire, enjoyment, worry, and frustration. Some of the examples he elaborates on involve higher-order beliefs. He asks (p. 48): "[H]ow can emotions help us explain behavior for which good explanations seem to be lacking?" We hope the framework we develop in this paper will be useful for providing answers.

# 8    Appendix

## 8.1    Extensive forms with observable actions

Fix a finite player set $N$ and finite action sets $A_i$ ($i \in N$). Let $A = \prod_{i \in N} A_i$. A history of length $\ell$ is a finite sequence of action profiles $h = (a^1, ..., a^\ell) \in A^\ell$. History $h = (a^1, ..., a^k)$ precedes history $\overline{h} = (\overline{a}^1, ..., \overline{a}^\ell)$, written $h \prec \overline{h}$, if $h$ is a prefix (initial subsequence) of $\overline{h}$, i.e. $k < \ell$ and $(a^1, ..., a^k) = (\overline{a}^1, ..., \overline{a}^k)$. In this case, we also write $\overline{h} = (h, \overline{a}^{k+1}, ..., \overline{a}^\ell)$. We let $\ell(h)$ denote the length of history $h$. The empty sequence (the history with zero length) is denoted by $h^0$. By convention $h^0$ precedes every proper history. A finite extensive form with

---

terial payoff functions are described. However, we have a sketch of proof regarding an equivalence of equilibrium predictions under the two approaches.

observable actions is a structure $\langle N, H \rangle$ where $H \subseteq \{h^0\} \cup \left( \bigcup_{\ell=1}^{L} A^\ell \right)$ is a finite set of histories with the following properties:[52]

- $h^0 \in H$.

- $\forall \overline{h} \in H$, if $h \prec \overline{h}$ then $h \in H$.

- $\forall h \in H$, $\{a \in A : (h, a) \in H\} = \prod_{i \in N} A_i(h)$ where

$$A_i(h) = \left\{ a_i \in A_i : \exists a_{-i} \in \prod_{j \neq i} A_j, (h, (a_i, a_{-i})) \in H \right\}$$

  is the set of possible actions of player $i$ at history $h$.

Note that $\langle H, \prec \rangle$ is a tree with distinguished root $h^0$; the symmetric closure of $\prec$ is denoted by $\preceq$.[53] We let $Z = \{h \in H : \prod_{i \in N} A_i(h) = \emptyset\}$ denote the set of terminal (or complete) histories.

We can now define the following derived elements:

- $S_i = \{s_i = (s_{i,h})_{h \in H} \in (A_i)^H : \forall h \in H \backslash Z, s_{i,h} \in A_i(h)\}$ is the set of strategies of player $i$, $S = \prod_{i=1}^{n} S_i$, $S_{-i} = \prod_{j \in N \backslash \{i\}} S_j$.

- $\zeta : S \to Z$ is the path function, that is, $z = (a^1, ..., a^K) = \zeta(s)$ iff $a^1 = (s_{i,h^0})_{i \in N}, \forall t \in \{1, ...K-1\}, a^{t+1} = (s_{i,(a^1,...,a^t)})_{i \in N}$.

- For any history $h \in H$, $S(h)$ is the set of strategy profiles consistent with history $h$, that is, $S(h) = \{s \in S : h \preceq \zeta(s)\}$. Since past actions are observed, it follows that $S(h) = \prod_{i=1}^{n} S_i(h)$, where $S_i(h)$ is the projection of $S(h)$ on $S_i$.

## 8.2 Proof of Theorem 9

We first show how to associate a consistent assessment $(\sigma, \beta(\sigma))$ to each behavioral profile $\sigma$. We define the first-order beliefs of $i$ corresponding to $\sigma$, $\beta^1(\sigma)$, as follows: $\beta_i^1(\sigma) = \mu_i^1$ where

$$\forall h \in H, \forall s_{-i} \in S_{-i}(h), \mu_i^1(s_{-i}|h) = \prod_{h' \, h} \prod_{j \neq i} \sigma_j(s_{j,h'}|h').$$

---

[52]Cf. Osborne & Rubinstein (1994, Chapter 6).
[53]Thus, $h \preceq h'$ iff either $h \prec h'$ or $h = h'$.

[Recall that we ignore $i$'s beliefs about himself; thus $\sigma_i$ occurs vacuously in the expression $\beta_i^1(\sigma)$.] It can be shown that $\mu_i^1 = \beta_i^1(\sigma)$ is a cps that satisfies the stochastic independence property, and furthermore that

$$\forall h \in H, \forall j \neq i, \forall a_j \in A_j(h), \ \mu_{ij}^1(S_j(h, a_j)|h) = \sigma_j(a_j|h).$$

Therefore the first-order beliefs obtained from $\sigma$ satisfy eq. (1) and conditions (a) and (b) in Definition 6 (consistency). The profile of belief hierarchies $\mu = \beta(\sigma)$ is obtained by condition (c) in Definition 6:

$$\begin{aligned} \forall i \ & \in \ N, \ \mu_i^1 = \beta_i^1(\sigma), \\ \forall i \ & \in \ N, \ \forall k > 1, \ \forall h \in H, \ \mu_i^k(\cdot|h) = \mu_i^{k-1}(\cdot|h) \times \delta_{t_{-i}^{k-1}} \ . \end{aligned}$$

Hence assessment $(\sigma, \beta(\sigma))$ is consistent. It is clear from the construction that $\beta(\cdot)$ is a continuous function.

**Definition 18** *Fix a strictly positive vector $\varepsilon = (\varepsilon_{i,h}(a_i)_{a_i \in A_i(h)})_{i \in N, h \in H \setminus Z}$ such that $\sum_{a_i \in A_i(h)} \varepsilon(a_i) < 1$ for all $h \in H \setminus Z$. An $\varepsilon$-psychological equilibrium is a behavioral strategy profile $\sigma$ such that $\forall i \in N$, $\forall h \in H$, $\forall a_i \in A_i(h)$, (i) $\sigma_i(a_i|h) \geq \varepsilon_{i,h}(a_i)$, (ii) $a_i \notin \arg\max_{a_i' \in A_i(h)} E_{\beta(\sigma)}[u_{i,h}|h, a_i']$ implies $\sigma_i(a_i|h) = \varepsilon_{i,h}(a_i)$.*

Let $\Sigma_\varepsilon$ denote the set of behavioral strategy profiles satisfying condition (i) of Definition 18 above and let $r_\varepsilon : \Sigma_\varepsilon \twoheadrightarrow \Sigma_\varepsilon$ denote the "$\varepsilon$-best response correspondence" that assigns to each profile $\sigma$ the subset of profiles in $\Sigma_\varepsilon$ satisfying condition (ii) of Definition 18, that is,

$$\begin{aligned} &r_{\varepsilon,i}(\sigma) \\ = \ & \{\sigma_i' \in \Sigma_{\varepsilon,i} : \forall h, \forall a_i, a_i \notin \arg\max_{a_i' \in A_i(h)} E_{\beta(\sigma)}[u_i|h, a_i'] \Rightarrow \sigma_i(a_i|h) = \varepsilon_{i,h}(a_i)\}, \end{aligned}$$

$$r_\varepsilon(\sigma) = \prod_{i \in N} r_{\varepsilon,i}(\sigma).$$

$r_{\varepsilon,i}(\sigma)$ is a nonempty convex subset of $\Delta(A_i(h))$. Since $E_{\beta(\sigma)}[u_i|h, a_i]$ is continuous in $(\sigma, \mu)$ and $\beta$ is a continuous function, $E_{\beta(\sigma)}[u_i|h, a_i]$ is continuous in $\sigma$. This implies that $r_{\varepsilon,i}(\sigma)$ has a closed graph. Thus, $r_\varepsilon(\cdot)$ is a nonempty convex valued correspondence with a closed graph from the compact and convex set $\prod_{h \in H \setminus Z} \Delta(A(h))$ to itself. By Kakutani theorem $r_\varepsilon(\cdot)$ has a fixed point, which is an $\varepsilon$-psychological equilibrium.

Fix a sequence $\varepsilon^k \to 0$ and a corresponding sequence of $\varepsilon^k$-psychological equilibria $\sigma^k$. By compactness, the sequence $(\sigma^k)$ has a limit point $\sigma^*$. We prove that $(\sigma^*, \beta(\sigma^*))$ is a psychological sequential equilibrium. Assessment $(\sigma^*, \beta(\sigma^*))$ is

consistent: to see this just note that, by continuity, $\beta(\sigma^*)$ is a limit point of $\beta(\sigma^k)$, and that the set of consistent assessments is closed. By continuity of $E_{(\sigma)}[u_i|h, a_i]$ in $\sigma$ (and finiteness of $A_i(h)$), for $k$ sufficiently large

$$\arg \max_{a_i \in A_i(h)} E_{(\sigma^*)}[u_i|h, a_i] = \arg \max_{a_i \in A_i(h)} E_{(\sigma^k)}[u_i|h, a_i].$$

This implies that

$$\mathrm{Supp}(\sigma_i^*(\cdot|h)) \subseteq \arg \max_{a_i \in A_i(h)} E_{(\sigma^*)}[u_i|h, a_i]$$

as required by Definition 8.∎

## 8.3  Results about interactive epistemology

We start with some preliminaries about rationality and backward induction on belief-induced decision trees, and then prove Lemmata 15, 16 and 17. Suppose that psychological payoff functions have the form $u_i : Z \times M \times S_{-i} \to \mathbb{R}$ [condition (4)]. Then, for any fixed hierarchy of cps' $\mu_i$, we obtain a well defined decision tree that can be solved by backward induction: define value functions $V_i : H \to \mathbb{R}$ and $\overline{V}_i : (H \backslash Z) \times A_i \to \mathbb{R} \cup \{-\infty\}$ as follows

- For terminal histories $z \in Z$, let

$$V_i(z) = \int_{S_{-i} \times \mathbf{M}_{-i}} u_i(z, s_{-i}, \boldsymbol{\mu}_{-i}, \boldsymbol{\mu}_i) f_{i,z}(\boldsymbol{\mu}_i)(ds_{-i}, d\boldsymbol{\mu}_{-i}).$$

- Assuming that $V_i(h, a)$ has been defined for the immediate successors $(h, a)$ of history $h$, let

$$\overline{V}_i(h, a_i) = \sum_{a_{-i} \in A_{-i}(h)} \mu_i^1(S_{-i}(h, a_{-i})|h) V_i(h, (a_i, a_{-i}));$$

for each $a_i \in A_i(h)$ and $V_i(h, a_i) = -\infty$ for $a_i \notin A_i(h)$;[54] then $V_i(h)$ is defined as

$$V_i(h) = \max_{a_i} \overline{V}_i(h, a_i)$$

(the maximum is well defined because $A_i$ is finite).

Recall that, for any strategy $s_i \in S_i$, $H_i(s_i) = \{h \in H \backslash Z : s_i \in S_i(h)\}$ denotes the set of histories allowed by $s_i$. The proof of the following result is available by request:

**Lemma 19** *Under condition (4) the sequential best reply correspondence* $r_i : M_i \twoheadrightarrow S_i$ *can be characterized as follows*

$$r_i(\boldsymbol{\mu}_i) = \left\{ s_i : \forall h \in H(s_i), s_{i,h} \in \arg \max_{a_i} \overline{V}_i(h, a_i) \right\}.$$

---

[54]We define $V_i(h, \cdot)$ outside $A_i(h)$ for notational convenience.

### 8.3.1 Proof of Lemma 15

By Lemma 19 $r_i(\mu_i) = \left\{ s_i : \forall h \in H(s_i), s_{i,h} \in \arg\max_{a_i} \overline{V}_{\cdot i}(h, a_i) \right\}$. Clearly, the RHS is nonempty. Therefore $r_i(\cdot)$ is nonempty-valued and $R_i$ is nonempty. The belief function $f_i$ is continuous (Lemma 3). If $u_i$ is also continuous, then $\mathrm{E}_{s_i, \cdot i}[u_i|h]$ is continuous (in $\mu_i$), which implies that $R_i$ is closed.∎

### 8.3.2 Proof of Lemma 16

We must show that for every closed event $E \in E_{-i}$, $\mathrm{SB}_i(E)$ is closed. $\mathrm{SB}_i(\emptyset) = \emptyset$, a closed set, by definition. Suppose that $E = \Omega_i \times E_{-i}$ where $E_{-i}$ is nonempty and closed. Recall that $\mathrm{SB}_i(E) = \bigcap_{h:[h] \cap E \neq \emptyset} \mathrm{B}_{i,h}(E)$. For each $h$,

$$\mathrm{B}_{i,h}(E) = S_i \times f_{i,h}^{-1}\left( \Delta(E_{-i} \cap (S_{-i}(h) \times \mathbf{M}_{-i})) \right) \times \Omega_{-i},$$

where for any measurable space $X$ and any $F \subseteq X$ we let $\Delta(F)$ denote the set of probability measures on $X$ that assign probability one to $F$. Note that if $F$ is closed, $\Delta(F)$ is also closed. The coordinate function $f_{i,h} : M_i \to \Delta(\Omega_{-i})$ is continuous and $M_{-i}$ is closed (Lemma 3); hence $E_{-i} \cap (S_{-i}(h) \times M_{-i})$, $\Delta(E_{-i} \cap (S_{-i}(h) \times M_{-i}))$ and $f_{i,h}^{-1}\left( \Delta(E_{-i} \cap (S_{-i}(h) \times \mathbf{M}_{-i})) \right)$ are closed. It follows that $\mathrm{B}_{i,h}(E)$ ($h \in H$) and $\mathrm{SB}_i(E)$ are closed.∎

### 8.3.3 Proof of Lemma 17

Let $\left\{ E^\ell \right\}_{\ell=0}^{\ell=k}$ be a decreasing sequence of nonempty events in $E$ ($\emptyset \neq E^k \subseteq E^{k-1} \subseteq ... \subseteq E^0$), we show that $\bigcap_{\ell=0}^{\ell=k} \mathrm{SB}(E^\ell)$ is also nonempty. For each $\ell$ and $i$, $E^\ell \in E$ can be written as $E^\ell = E_i^\ell \times E_{-i}^\ell$, where $E_{-i}^\ell \subseteq \Omega_{-i}$, and by definition of $\mathrm{SB}(\cdot)$

$$\bigcap_{\ell=0}^{\ell=k} \mathrm{SB}(E^\ell) = \bigcap_{i \in N} \bigcap_{\ell=0}^{\ell=k} \mathrm{SB}_i(\Omega_i \times E_{-i}^\ell).$$

Therefore we must show that $\bigcap_{\ell=0}^{\ell=k} \mathrm{SB}_i(\Omega_i \times E_{-i}^\ell) \neq \emptyset$ ($i \in N$). Let $\Delta^H(\Omega_{-i}; E_{-i}^\ell)$ denote the set of cps' $\mu \in \Delta^H(\Omega_{-i})$ such that $\mu(E_{-i}^\ell|h) = 1$ for each $h$ such that $E_{-i}^\ell \cap (S_{-i}(h) \times M_{-i}) \neq \emptyset$. Note that

$$\bigcap_{\ell=0}^{\ell=k} \mathrm{SB}_i(\Omega_i \times E_{-i}^\ell) = S_i \times f_i^{-1}\left( \bigcap_{\ell=0}^{\ell=k} \Delta^H(\Omega_{-i}; E_{-i}^\ell) \right) \times \Omega_{-i}.$$

We show below that $\bigcap_{\ell=0}^{\ell=k} \Delta^H(\Omega_{-i}; E_{-i}^\ell) \neq \emptyset$. Since $f_i$ is onto (Lemma 3), it follows that $f_i^{-1}\left( \bigcap_{\ell=0}^{\ell=k} \Delta^H(\Omega_{-i}; E_{-i}^\ell) \right) \neq \emptyset$. Hence $\bigcap_{\ell=0}^{\ell=k} \mathrm{SB}_i(\Omega_i \times E_{-i}^\ell) \neq \emptyset$.

We show that $\bigcap_{\ell=0}^{\ell=k} \Delta^H(\Omega_{-i}; E_{-i}^\ell) \neq \emptyset$ with a recursive construction. Say that $h$ is 'reached' by probability measure $\nu \in \Delta(\Omega_{-i})$ if $\nu(S_{-i}(h) \times M_{-i}) > 0$. Note that if $h$ is reached by $\nu$, every predecessor of $h$ is also reached by $\nu$. Say that $\mu(\cdot|h)$ is 'derived' from $\nu$, where $\nu$ reaches $h$, if for every Borel set $F_{-i} \subseteq \Omega_{-i}$

$$\mu(F_{-i}|h) = \frac{\nu(F_{-i} \cap (S_{-i}(h) \times \mathbf{M}_{-i}))}{\nu(S_{-i}(h) \times \mathbf{M}_{-i})}.$$

Pick any probability measure $\nu$ in the (nonempty) set $\Delta(E_{-i}^k)$. For each $h$ reached by $\nu$ let $\mu(\cdot|h)$ be derived from $\nu$. Thus, $\mu(\cdot|h)$ has been defined for a nonempty set of histories closed w.r.t. precedence (that is, if $h$ is in the set every predecessor of $h$ is in the set), the set is nonempty because it contains the initial history $h^0$. Now suppose that $\mu(\cdot|h)$ has been defined for some set of histories $\hat{H}$ closed w.r.t. precedence. If $\hat{H} \neq H$, for each $h \in H \backslash \hat{H}$ such that the immediate predecessor of $h$ belongs to $\hat{H}$, pick a probability measure $\nu_h$ in the set $\Delta(E_{-i}^{\ell(h)} \cap (S_{-i}(h) \times M_{-i}))$, where $\ell(h)$ is the highest index $\ell \in \{-1, 0, ..., k\}$ such that $E_{-i}^\ell \cap (S_{-i}(h) \times M_{-i}) \neq \emptyset$, and by convention we let $E^{-1} = \Omega_{-i}$. Let $\mu(\cdot|h')$ be derived from $\nu_h$ whenever $h'$ weakly follows $h$ and is reached by $\nu_h$. Now $\mu(\cdot|h)$ is defined for a set of histories $\hat{H}'$ closed under the precedence relation and strictly larger than $\hat{H}$. Proceed in this way until the whole $H$ is covered. We claim that the resulting vector of probability measures $(\mu(\cdot|h))_{h \in H}$ is a cps $\mu \in \bigcap_{\ell=0}^{\ell=k} \Delta^H(\Omega_{-i}; E_{-i}^\ell)$.

To see that $(\mu(\cdot|h))_{h \in H}$ is a cps we only have to check that the 'chain rule' (3) in Definition 1. Suppose that $h$ precedes $h'$. To write formulas more transparently, let $C = S_{-i}(h) \times M_{-i}$, $C'_{-i} = S_{-i}(h') \times M_{-i}$, $\mu(\cdot|h) = \mu(\cdot|C_{-i})$, $\mu(\cdot|h') = \mu(\cdot|C'_{-i})$. Since $h$ precedes $h'$, $S_{-i}(h') \subseteq S_{-i}(h)$, hence $C'_{-i} \subseteq C_{-i}$. If $h'$ is not reached by $\mu(\cdot|C_{-i})$ then (3) holds trivially as $0 = 0$. If $h'$ is reached by $\mu(\cdot|C_{-i})$, then $\mu(\cdot|C_{-i})$ and $\mu(\cdot|C'_{-i})$ are both derived from the same measure – say $\nu \in \Delta(\Omega_{-i})$ – reaching both $h$ and $h'$; thus, for every Borel set $F_{-i} \subseteq C'_{-i}$

$$\mu(F_{-i}|C_{-i}) = \frac{\nu(F_{-i})}{\nu(C_{-i})} = \frac{\nu(F_{-i})}{\nu(C'_{-i})} \frac{\nu(C'_{-i})}{\nu(C_{-i})} = \mu(F_{-i}|C'_{-i})\mu(C'_{-i}|C_{-i}).$$

To see that $\mu \in \bigcap_{\ell=0}^{\ell=k} \Delta^H(\Omega_{-i}; E_{-i}^\ell)$, note that by construction $\mu(E^{\ell(h)}|h) = 1$ for all $h \in H$. Suppose that, for any index $\ell \in \{0, ..., k\}$ and any $h \in H$, $E_{-i}^\ell \cap (S_{-i}(h) \times M_{-i}) \neq \emptyset$. Then $\ell(h) \geq \ell$ and $\mu(E^\ell|h) \geq \mu(E^{\ell(h)}|h) = 1$; hence $\mu(E^\ell|h) = 1$ as desired.∎

# References

[1] ALIPRANTIS, C. and K. BORDER (1999): Infinite Dimensional Analysis. Berlin: Springer-Verlag (2nd edition).

[2] AUMANN, R.J. and A. BRANDENBURGER (1995): "Epistemic Conditions for Nash Equilibrium," Econometrica, 63, 1161-1180.

[3] BACHARACH, M., G. GUERRA and D.J. ZIZZO (2001): "Is Trust Self-Fulfilling? An Experimental Study," mimeo.

[4] BATTIGALLI, P. (1996): "Strategic Independence and Perfect Bayesian Equilibria," Journal of Economic Theory, 70, 201-234.

[5] BATTIGALLI, P. and M. SINISCALCHI (1999): "Hierarchies of Conditional Beliefs and Interactive Epistemology in Dynamic Games," Journal of Economic Theory, 88, 188-230.

[6] BATTIGALLI, P. and M. SINISCALCHI (2002): "Strong Belief and Forward Induction Reasoning," Journal of Economic Theory, 106, 356-391.

[7] BELL, D.E. (1982): "Regret in Decision Making under Uncertainty," Operations Research, 30, 961-981.

[8] BEN PORATH, E. and E. DEKEL (1992): "Signaling Future Actions and the Potential for Sacrifice, " Journal of Economic Theory, 57, 36-51.

[9] BERNHEIM, D. (1994): "A Theory of Conformity," Journal of Political Economy, 102, 841-77.

[10] BRANDENBURGER, A. and E. DEKEL (1993): "Hierarchies of Beliefs and Common Knowledge," Journal of Economic Theory, 59, 189-198.

[11] CHARNESS, G. and M. DUFWENBERG (2004): "Promises and Partnership", mimeo, UC Santa Barbara and University of Arizona.

[12] DUFWENBERG, M. (1995): "Time-Consistent Wedlock with Endogenous Trust", in On Rationality and Belief Formation in Games, Doctoral Dissertation, Department of Economics, Uppsala University.

[13] DUFWENBERG, M. (2002): "Marital Investment, Time Consistency and Emotions," Journal of Economic Behavior and Organization, 48, 57-69.

[14] DUFWENBERG, M. and U. GNEEZY (2000): "Measuring beliefs in an experimental lost wallet game," Games and Economic Behavior, 30, 163-182.

[15] DUFWENBERG, M. and G. KIRCHSTEIGER (2004): "A Theory of Sequential Reciprocity," Games and Economic Behavior, 47, 268-298.

[16] DUFWENBERG, M. and M. LUNDHOLM (2001): "Social Norms and Moral Hazard," Economic Journal, 111, 506-525.

[17] ELSTER, J. (1998): "Emotions and Economic Theory," Journal of Economic Literature, 36, 4774.

[18] FALK, A. and U. FISCHBACHER (1998): "A Theory of Reciprocity," mimeo.

[19] FEHR, E. and S. GÄCHTER (2000): "Fairness and Retaliation: The Economics of Reciprocity," Journal of Economic Perspectives, 14, 159-181.

[20] FUDENBERG, D. and D.K. LEVINE (1993): "Self-Confirming Equilibrium," Econometrica, 61, 523-545.

[21] FUDENBERG, D. and J. TIROLE (1991a): Game Theory. Cambridge MA: MIT Press.

[22] FUDENBERG, D. and J. TIROLE (1991b): "Perfect Bayesian Equilibria," Journal of Economic Theory, 53, 236-260.

[23] GEANAKOPLOS, J., D. PEARCE and E. STACCHETTI (1989): "Psychological Games and Sequential Rationality," Games and Economic Behavior, 1, 60-79.

[24] GEANAKOPLOS, J. (1996): "The Hangman Paradox and the Newcomb's Paradox as Psychological Games," Cowles Foundation Discussion Paper No. 1128.

[25] GILBOA, I. and D. SCHMEIDLER (1988): "Information Dependent Games: Can Common Sense Be Common Knowledge?" Economics Letters, 1988, 27, 215-221.

[26] GUERRA, G. and D.J. ZIZZO (2004): "Trust Responsiveness and Beliefs," Journal of Economic Behavior and Organization, 55, 25-30.

[27] GUL, F. and W. PESENDORFER (2004): "The Canonical Type Space for Interdependent Preferences," mimeo, Princeton University.

[28] HARSANYI, J. (1967-68): "Games of Incomplete Information Played by Bayesian Players. Parts I, II, III," Management Science, 14, 159-182, 320-334, 486-502.

[29] HUANG, P.H. and WU, H.-M. (1994): "More Order without More Law: A Theory of Social Norms and Organizational Cultures," Journal of Law, Economics and Organization, 12, 390406.

[30] HUCK, S. and D. KÜBLER (2000): "Social Pressure, Uncertainty, and Co-operation," Economics of Governance, 1, 199-212.

[31] KECHRIS, A. (1995): Classical Descriptive Set Theory, Berlin: Sringer Verlag.

[32] KOHLBERG E. and P. RENY (1997): "Independence on Relative Probability Spaces and Consistent Assessments in Game Trees," Journal of Economic Theory, 75, 280-313.

[33] KOLPIN, V. (1992): "Equilibrium Refinements in Psychological Games," Games and Economic Behavior, 4, 218-231.

[34] KROHNE, H.W. (2003): "Individual differences in emotional reactions and coping," in R. J. Davidson, K. R. Scherer & H. H. Goldsmith (Eds.), Handbook of Affective Sciences, pp. 698-725. New York: Oxford University Press.

[35] KUHN, H.W. (1953): "Extensive Games and the Problem of Information," in Contributions to the Theory of Games II, ed. by H. W. Kuhn and A. W. Tucker. Princeton: Princeton University Press, pp. 193-216.

[36] LI, J. (2005): "The Power of Convention: A Theory of Social Preferences," mimeo, University of Pennsylvania.

[37] LOOMES, G. and R. SUGDEN (1982), "Regret Theory: An Alternative Theory of Rational Choice under Uncertainty," Economic Journal, 92, 805-824.

[38] MERTENS J.F. and S. ZAMIR (1985): "Formulation of Bayesian Analysis for Games with Incomplete Information," International Journal of Game Theory, 14, 1-29.

[39] OSBORNE, M. and A. RUBINSTEIN (1994): A Course in Game Theory. Cambridge MA: MIT Press.

[40] PEARCE, D. (1984): "Rationalizable Strategic Behavior and the Problem of Perfection," Econometrica, 52, 1029-1050.

[41] RABIN, M. (1993): "Incorporating Fairness into Game Theory and Economics," American Economic Review, 83, 1281-1302.

49

[42] RENY, P. (1992): "Backward Induction, Normal Form Perfection and Explicable Equilibria," Econometrica, 60, 626-649.

[43] RÊNYI, A. (1955): "On a New Axiomatic Theory of Probability," Acta Mathematica Academiae Scientiarum Hungaricae, 6, 285-335.

[44] RUBINSTEIN, A. (1991): "Comments on the Interpretation of Game Theory," Econometrica, 59, 909-904.

[45] RUFFLE, B.J. (1999): "Gift Giving with Emotions," Journal of Economic Behavior and Organization, 39, 399-420.

[46] SEGAL, U. and J. SOBEL (2003): "Tit for Tat: Foundations for Preferences for Reciprocity in Strategic Settings," mimeo, UCLA and Boston College.

[47] TANGNEY, J.P. (1995): "Recent Advances in the Empirical Study of Shame and Guilt," American Behavioral Scientist, 38, 11321145.

[48] TAN, T. and S. WERLANG (1988): "The Bayesian Foundation of Solution Concepts of Games," Journal of Economic Theory, 45, 370-391.