



# Epistemic game theory without types structures: An application to psychological games <sup>☆</sup>

Pierpaolo Battigalli <sup>a,\*</sup>, Roberto Corrao <sup>b</sup>, Federico Sanna <sup>c</sup>

<sup>a</sup> Department of Economics and IGIER, Bocconi University, Via Roberto Sarfatti 25, 20136 Milan, Italy

<sup>b</sup> Department of Economics, MIT, Memorial Drive 50, 02142 Cambridge (MA), USA

<sup>c</sup> Lyra Partners, Via Broletto 35, 20121 Milano (MI), Italy

## ARTICLE INFO

### Article history:

Received 15 January 2019

Available online 31 December 2019

### JEL classification:

C72

C73

D82

### Keywords:

Epistemic game theory

Hierarchies of beliefs

Consistency

Subjective rationality

Strong rationalizability

Psychological games

## ABSTRACT

We consider multi-stage games with incomplete information, and we analyze strategic reasoning by means of epistemic events within a “total” state space made of all the profiles of behaviors (paths of play) and possibly incoherent infinite hierarchies of conditional beliefs. Thus, we do not rely on types structures, or similar epistemic models. Subjective rationality is defined by the conjunction of coherence of belief hierarchies, rational planning, and consistency between plan and on-path behavior. Since consistent hierarchies uniquely induce beliefs about behavior and belief hierarchies of others, we can define rationality and common strong belief in rationality, and analyze their behavioral and low-order beliefs implications, which are characterized by strong rationalizability. Our approach allows to extend known techniques to the epistemic analysis of psychological games where the utilities of outcomes depend on beliefs of order  $k$  or lower. This covers almost all applications of psychological game theory.

© 2019 Elsevier Inc. All rights reserved.

## 1. Introduction

**Epistemic game theory** (EGT) is the formal analysis of players' interactive strategic reasoning in games.<sup>1</sup> Such analysis posits, or obtains by construction, a set **states of the world**  $\Omega$  so that each  $\omega \in \Omega$  is an all-encompassing implicit or explicit specification of all the relevant aspects of the strategic situation, including what players do and how they think about each other's behavior and beliefs. This permits the definition of events (measurable subsets of  $\Omega$ ) such as “players are rational” [viz.  $R \subseteq \Omega$ ] and “it is common belief that players are rational” [viz.  $CB(R) \subseteq \Omega$ ]. These events are the **epistemic assumptions** of interest and relate behavior to beliefs. The typical theorem of epistemic game theory provides a characterization of the interesting implications of such epistemic assumptions, such as the behavioral implications. For example, given appro-

<sup>☆</sup> We have benefited from helpful comments from two anonymous referees, Nicodemo De Vito, David Ruiz Gomez, Giacomo Lanzani, Muhamet Yildiz, and the seminar participants at the 3rd Annual Workshop on Behavioural Game Theory in Norwich, 2017, and at LOFT, in Milan, 2018. We thank Federico Bobbio, Carlo Cusumano, Francesco Fabbri, Davide Ferri, and Giulio Principi for outstanding research assistantship. Financial support of ERC (grant 324219) and of the Marco Fanno scholarship are gratefully acknowledged.

\* Corresponding author.

E-mail address: [pierpaolo.battigalli@unibocconi.it](mailto:pierpaolo.battigalli@unibocconi.it) (P. Battigalli).

<sup>1</sup> See, for example, the survey by Dekel and Siniscalchi (2015), or the textbook by Perea (2012), and the relevant references therein.

appropriate definitions of “rationality,” “common belief,” and “rationalizability,” we have the following theorem<sup>2</sup>: *Players’ behavior is consistent with  $R \cap CB(R)$  if and only if it is rationalizable.*

**Standard approach** The standard approach of EGT is to posit or construct  $\Omega$  so that it is true at every state  $\omega$  in  $\Omega$  that every player is **cognitively rational**, that is, that his system of beliefs satisfies appropriate coherence requirements. Specifically, each state  $\omega$  determines for each player  $i$  a hierarchy of beliefs such that the “first-order” belief of  $i$  about the behavior of other players  $-i$  is the marginal of “higher-order” joint beliefs about the behavior and beliefs of  $-i$ , and similar coherence restrictions hold for higher-order beliefs. Since cognitive rationality holds at every state, it has to be commonly believed at every state that players are cognitively rational. Say that an event is **transparent** if it is true and commonly believed. With this, *the standard approach assumes that cognitive rationality is transparent.*

Yet, rationality has also a behavioral aspect: roughly, a player is **rational** at  $\omega$  if he is cognitively rational at  $\omega$  and his behavior at  $\omega$  is a best reply to his beliefs at  $\omega$ . Most contributions to EGT do *not* assume that players are rational at every state. In particular, allowing for states where players are irrational is important in the epistemic analysis of multi-stage games that contain histories inconsistent with either mere rationality, or with rationality and common belief in rationality. Indeed, assuming that rationality holds at every state would imply “by fiat” that the actions in such histories cannot be chosen even if they comply with the rules of the game. Several authors (including ourselves) find this feature conceptually problematic. Thus, standard EGT postulates transparency of cognitive rationality, but instead treats rationality and common belief in rationality as a property that holds only in some states. We accept the latter, but question the former: *Why should cognitive rationality be transparent while rationality holds only at some states?*

In our view, the reason is more technical than conceptual: the standard approach allows to work with epistemic structures. In particular, much of the EGT literature on multi-stage games works with type structures, whereby each player is uncertain about the coplayers’ behavior and a type of a player corresponds to a hierarchy of beliefs satisfying coherence and common belief in coherence. It is argued that this is without loss of generality, because one can show that the space of profiles of behaviors and belief hierarchies satisfying transparency of coherence gives the “largest” type structure.<sup>3</sup>

**Our approach** Working with type structures is traditional and in many ways convenient, but we show in this paper that eschewing them is both possible and fruitful. *We regard cognitive rationality simply as an aspect of rationality that—like rationality itself—holds only in some states.* Hence, we consider the “total” space of all profiles of behaviors and beliefs, including incoherent beliefs. Rationality of player  $i$  corresponds to the set of states where  $i$  has coherent beliefs and his behavior is a sequential best reply to such beliefs. Since the belief-hierarchy of a rational player is coherent, it induces a belief (a conditional probability system) about behaviors and beliefs hierarchies of the opponents.<sup>4</sup> Leveraging on this, we prove that, in multi-stage games with possibly incomplete information, strong rationalizability characterizes the behavioral implications of rationality and common strong belief in rationality, which represents forward-induction reasoning.

We take advantage of our “fresh start” to introduce other innovations compared to standard epistemic game theory. Indeed, we take as primitive players’ *uncertainty about the path of play*, rather than uncertainty about contingent behavior. With this, the only mathematical objects that look like “strategies” are players’ (marginal) systems of beliefs about their own actions conditional of reaching non-terminal histories. These are mere plans. Rationality is modeled as coherence of beliefs, rational planning (incentive compatibility of plans), and on-path consistency of plan and actions.<sup>5</sup>

**Application to psychological games** Besides its conceptual appeal, our approach has also a technical advantage: it allows us to extend to so called “psychological games” the techniques used by Battigalli and Tebaldi (2019) in the epistemic analysis of a class of infinite dynamic games. In a **psychological game** utilities of outcomes depend on beliefs. This allows to capture a wide range of emotional or otherwise psychological aspects of choice.<sup>6</sup> In applications, the utility of outcomes is assumed to depend only on beliefs of the first  $k$  orders, e.g., only on the first-order beliefs of everybody. Dependence on beliefs up to order  $k$  allows for a tractable definition of rationalizability by means of iterated elimination of non-best replies to beliefs of order  $k + 1$ . The technical problem is to show that a  $k + 1$ th-order belief that justifies a player’s behavior as rationalizable can be extended to an infinite hierarchy of beliefs that makes such behavior consistent with rationality and common strong belief in rationality. Battigalli and Tebaldi (2019) rely on the possibility to *factorize* the uncertainty space of any player  $i$  as  $\Omega_{-i} = U_{-i} \times M_{-i}$ , where  $U_{-i}$  is a set of utility-relevant states and  $M_{-i}$  is a space of opponents’ beliefs (probability measures). In a psychological game,  $U_{-i}$  should be the set of possible behaviors and opponents’ beliefs up to order  $k$ , and  $M_{-i}$  should be the set of opponents’ beliefs of order  $k + 1$  or higher. Yet, when working with type structures,  $\Omega_{-i}$  cannot be factorized in this way, because low-order beliefs must be the marginal of higher-order beliefs; thus,  $\Omega_{-i} \subsetneq U_{-i} \times M_{-i}$ , which prevents the application of the aforementioned methods. We instead work with the space of *all* belief hierarchies,

<sup>2</sup> One can give versions of this result for different decision criteria (forms of rationality), different kinds of game (e.g., simultaneous, or sequential, with complete or incomplete information), and—correspondingly—different definitions of rationalizability.

<sup>3</sup> See Battigalli and Siniscalchi (1999).

<sup>4</sup> See Proposition 1 in Battigalli and Siniscalchi (1999) and Proposition 1 in Brandenburger and Dekel (1993).

<sup>5</sup> For more on this, see the discussion in Section 8.

<sup>6</sup> See the survey by Battigalli and Dufwenberg (2019).

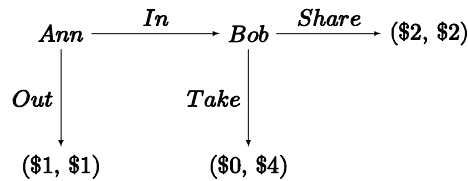


Fig. 1. Trust Game form.

including the incoherent ones, which implies that the relevant factorization holds. With this, we prove the results stated above for multi-stage games where psychological utility depends on beliefs of some given order  $k$ , or lower. To be clear, we are not claiming that our approach is necessary for an epistemic analysis of psychological games,<sup>7</sup> we only argue that our approach has an independent conceptual motivation and it also has the advantage of allowing the extension of known methods from standard games to psychological games.

**Heuristic example** We illustrate the main ideas of this paper through a simple example. We model guilt aversion in a two-person game form with monetary payoffs<sup>8</sup> as the desire—other things being equal—not to disappoint the other player. The disappointment of player  $j$  at terminal history  $z$  given first-order belief  $\mu_{j,1}$  is the difference, if positive, between her or his (initially) expected monetary payoff and his realized monetary payoff at  $z$ .<sup>9</sup> Thus, the psychological utility of player  $i \neq j$  depends on  $\mu_{j,1}$  (unknown to  $i$ ) as in the following equation

$$u_i(z, \theta_i, \mu_{j,1}) = \pi_i(z) - \theta_i \max\{0, \mathbb{E}_{\mu_{j,1}}(\pi_j) - \pi_j(z)\}, \quad \underline{\theta}_i \leq \theta_i \leq \bar{\theta}_i, \quad (1)$$

where  $\pi_i(z)$  denotes the monetary payoff of  $i$  at  $z$ , and  $\theta_i$  is a sensitivity parameter known to  $i$ . Furthermore, it is common knowledge that  $\theta_i$  belongs to the compact interval  $\Theta_i = [\underline{\theta}_i, \bar{\theta}_i]$  with  $\underline{\theta}_i \geq 0$ . Note that  $u_i$  is a kind of state-dependent utility function, because it depends on a feature of  $j$ , his first-order belief, that  $i$  neither knows nor controls. To assess the expected utility of his actions, player  $i$  has to consult his second-order beliefs  $\mu_{i,2}$ , that is, his beliefs about behavior and the personal features of  $j$ , whether exogenous ( $\theta_j$ ) or affected by strategic thinking ( $\mu_{j,1}$ ). Suppose that the game form with monetary payoffs is the Trust Game depicted in Fig. 1. Adding psychological utility functions as in Eq. (1) we obtain a first-order psychological game. We provide an informal analysis based on the epistemic assumptions of (correct) common strong belief in rationality and the solution concept that characterizes its utility-relevant implications, strong rationalizability (cf. Battigalli and Siniscalchi, 2002).

To simplify the analysis we assume throughout that Ann is commonly known to be selfish (and risk neutral), i.e.,  $\underline{\theta}_a = \bar{\theta}_a$ . Since Ann can secure payoff \$1, if Bob believes whenever possible (=strongly believes) in her rationality, then he would infer from *In* that her disappointment after *Take* would be at least 1. Hence, by Eq. (1), the expected utility of *Take* given *In* satisfies  $\mathbb{E}_{\mu_{b,2}}(u_b|In, Take) \leq 4 - \theta_b$ . Since Ann cannot be disappointed by *Share*, Bob's utility of *Share* given *In* is 2, and Bob would certainly *Share* if  $4 - \theta_b < 2$ , or  $\theta_b > 2$ .

Let us first analyze the case in which Bob is commonly known to be “sufficiently guilt averse”:  $\underline{\theta}_b > 2$ . If Ann (on top of being rational) believes that Bob is rational and that he strongly believes in her rationality, then she expects to get \$2 if she goes *In* and so she does. If Bob believes all this, he indeed expects *In*. Theorem 1 shows that strong rationalizability, an iterated elimination procedure, characterizes the utility-relevant implications of rationality, (correct) strong belief in rationality, and so on. **Step 1** of the procedure eliminates all path-belief pairs  $(z, \mu_{a,1})$  such that  $z = (In, a_b)$ , and  $\mathbb{E}_{\mu_{a,1}}(\pi_a|In) < 1$  or  $\mu_{a,1}(In) = 0$ . Note that  $((In, a_b), \mu_{a,1})$  with  $\mu_{a,1}(In) = 0$  is eliminated because  $\mu_{a,1}(In) = 0$  means that Ann plans to go *Out*, but at this personal state she does the opposite, and such inconsistency between plan and behavior is a form of irrationality.<sup>10</sup> Since  $\underline{\theta}_b > 2$  and all the non-eliminated pairs  $((In, a_b), \mu_{a,1})$  satisfy  $\mathbb{E}_{\mu_{a,1}}(\pi_a|In) \geq 1$ , **step 2** eliminates all  $(z, \mu_{b,1})$  such that either  $z = (In, Take)$  or  $\mu_{b,1}(Share|In) < 1$  (or both), because there is no second-order belief  $\mu_{b,2}$  assigning probability 1 to  $\{\mu_{a,1} : \mathbb{E}_{\mu_{a,1}}(\pi_a|In) \geq 1\}$  and such that *Take* is a best reply to  $\mu_{b,2}(\cdot|In)$ . Note, the procedure eliminates paths of play (behaviors) and first-order beliefs, where the latter include players' plans. In this case Bob must plan to *Share* if given the opportunity and must actually *Share* if Ann goes *In*. Thus, we look at actual behavior (paths of play) and plans separately, and we require that rational players (1) plan to choose best replies to their beliefs (rational planning) and never take actions that they planned not to take (material consistency). Steps  $n > 1$ , at any given history/node  $h$ , take into account only non-eliminated pairs, as long as they allow for  $h$ , which captures strong belief in previous steps.

Case  $\bar{\theta}_b < 1$ , in which Bob is commonly known to be “sufficiently selfish” is easier. The first step for Ann is the same as before, but now mere rationality has implications for Bob as well: even if, upon observing *In*, he were certain that Ann

<sup>7</sup> See our comments on the literature.

<sup>8</sup> Note, the game form represents only the rules of the game, not players' preferences. Here we assume that the consequences of players' behavior are monetary payoffs, which do not necessarily represent players' preferences.

<sup>9</sup> See Battigalli and Dufwenberg (2019), Battigalli et al. (2019a), and the relevant references therein.

<sup>10</sup> Similarly, all pairs  $(z, \mu_{a,1})$  such that  $z = (Out)$ , and either  $\mathbb{E}_{\mu_{a,1}}(\pi_a|In) > 1$  or  $\mu_{a,1}(Out) = 0$  are eliminated as well.

expected to get \$2, he would still *Take*, because his guilt sensitivity is too low. In this case, only *Out* is consistent with rationality and belief in rationality.

If  $\theta_b < 2$  and  $\theta_b > 1$ , the analysis is more complex: for Bob, we have to look at triplets  $(z, \theta_b, \mu_{b,1})$  and delete them if  $z = (In, a_b)$  and  $a_b$  is not a best reply for “utility type”  $\theta_b$  to any  $\mu_{b,2}(\cdot|In)$  satisfying the belief restrictions implied by previous steps (or if  $\mu_{b,1}(a_b|In) = 0$ ). Anyway, the upshot is that in this case information incompleteness prevents definite predictions: every path is consistent with rationality and common strong belief in rationality.

**Related literature** To the best of our knowledge, this is the first paper on epistemic game theory that does not assume transparency of coherence and hence does not rely on standard epistemic structures. As we explained above, our epistemic justification of solution concepts adapts techniques borrowed from Battigalli and Tebaldi (2019). Our analysis of hierarchies of conditional beliefs builds on Battigalli and Siniscalchi (1999), and their result (Proposition 1) that coherent hierarchies of conditional beliefs can be homeomorphically mapped to conditional beliefs about external uncertainty and (coherent as well incoherent) belief hierarchies of the other players.<sup>11</sup> Our analysis of common strong belief in rationality (forward-induction reasoning) builds on Battigalli and Siniscalchi (2002) and Battigalli et al. (2013). In the latter paper, as in ours, primitive uncertainty concerns the path of play. Psychological games were first analyzed by Geanakoplos et al. (1989). Here we borrow from and modify the extended framework of Battigalli and Dufwenberg (2009), who also provide the first epistemic analysis of (strong) rationalizability in psychological games. The most important difference with Battigalli and Dufwenberg (2009) is that they do not assume a finite upper bound on the order of beliefs that affect psychological utility; hence, they cannot give a tractable definition of (strong) rationalizability as iterated elimination of non-best replies. Also, Battigalli and Dufwenberg (2009) assume transparency of coherence and of consistency between plan and behavior, and that psychological utility does not depend on one’s own plan, while we dispense with these assumptions. In particular, this allows for an explicit analysis of players’ inferences about coplayers’ intentions and it introduces the possibility of dynamically inconsistent preferences. Both features are crucial in some important applications of psychological game theory.<sup>12</sup> Finally, Jagau and Perea (2017, 2018) analyze rationality and common belief in rationality in simultaneous-move games where psychological utility depends only on *initial* beliefs. Our analysis differs from theirs in several aspects. First, we consider multi-stage games<sup>13</sup> where psychological utility may depend on initial as well as updated beliefs, including *terminal* beliefs. Therefore, we are able to model—both in one-stage and in multi-stage games—backward and forward-induction reasoning, as well as concerns for the opinions of others (see Battigalli and Dufwenberg, 2009, 2019). Second, Jagau and Perea (2017) rely on the standard type-structure approach; hence, they use different methods.

**Outline** The rest of the paper is organized as follows. Section 2 reviews some mathematical preliminaries. Section 3 presents multi-stage game forms. Section 4 defines hierarchies of beliefs and the total state space of the game. Section 5 defines  $k$ -th order psychological games and subjective rationality. Section 6 defines strong rationalizability. Section 7 states and proves our main result, Theorem 1, which provides an epistemic justification of strong rationalizability. Finally, Section 8 discusses possible extensions of our work.

**2. Mathematical preliminaries**

For every set  $X$  and for every  $n \in \mathbb{N}$ , let  $X^n$  denote the  $n$ -fold product of the set  $X$  with generic element denoted by  $x^n$ . By convention, we let  $X^0 = \{\emptyset\}$ , where  $\emptyset$  denotes the empty sequence. Fix a (non-ordered) index set  $I$ , a profile of sets  $(X_i)_{i \in I}$ , and let  $\prod_{i \in I} X_i$  denote the set of all selections from correspondence  $i \mapsto X_i$ , that is, the set of all functions  $f : I \rightarrow \bigcup_{i \in I} X_i$  such that  $f(i) \in X_i$ . In other words, the Cartesian product  $\prod_{i \in I} X_i$  is regarded as a set of functions and the order of “factors” is irrelevant.

Given any compact metrizable topological space  $\Omega$ , let  $\mathcal{B}(\Omega)$  denote its Borel sigma-algebra and let  $\Delta(\Omega)$  denote the space of all the Borel probability measures over  $\mathcal{B}(\Omega)$ . We always endow  $\Delta(\Omega)$  with the topology of weak convergence of probability measures. It follows that  $\Delta(\Omega)$  is a compact metrizable topological space (see for example Aliprantis and Border, 2006, Chapter 15). We always endow finite spaces with the discrete topology, product of topological spaces with the product topology and subsets of topological spaces with the relative topology. Given a measurable subset  $F$  of  $\Omega$ , let  $\mathcal{B}(\Omega) \cap F$  denote the relative Borel sigma-algebra of  $F$ .

**Definition 1.** Let  $\Omega$  be compact metrizable and let  $\mathcal{F}$  be a countable collection of *clopen* subsets of  $\Omega$ . The pair  $(\Omega, \mathcal{F})$  is called **conditional measurable space**.

The following definition is key in our analysis.

**Definition 2.** Let  $(\Omega, \mathcal{F})$  be a conditional measurable space. A **conditional probability system (CPS)** over  $(\Omega, \mathcal{F})$  is an array  $\mu \in [\Delta(\Omega)]^{\mathcal{F}}$  that satisfies:

<sup>11</sup> This is an extension of Proposition 1 in Brandenburger and Dekel (1993), which concerns hierarchies of probability measures.  
<sup>12</sup> For a comprehensive analysis of the main applications of psychological game theory see Battigalli and Dufwenberg (2019) and Battigalli et al. (2019a).  
<sup>13</sup> Including simultaneous games as a special case.

- (Knowledge implies belief) For all  $F \in \mathcal{F}$  and for all  $E \in \mathcal{B}(\Omega)$

$$F \subseteq E \Rightarrow \mu(E|F) = 1.$$

- (Chain Rule) For all  $F_1, F_2 \in \mathcal{F}$  with  $F_2 \subseteq F_1$  and for all  $E \in \mathcal{B}(\Omega) \cap F_2$

$$\mu(E|F_1) = \mu(E|F_2) \mu(F_2|F_1).$$

Let  $\Delta^{\mathcal{F}}(\Omega) \subseteq [\Delta(\Omega)]^{\mathcal{F}}$  denote the set of CPSs on  $(\Omega, \mathcal{F})$ .

The defining properties of CPSs essentially say that the rules of conditional probability apply whenever possible. In particular, they imply that, for each  $F \in \mathcal{F}$ , each measurable partition  $\{F_1, \dots, F_k\}$  of  $F$  and each measurable subset  $E$  of  $F$ ,

$$\mu(E|F) = \sum_{\ell=1}^k \mu(E|F_{\ell}) \mu(F_{\ell}|F).$$

Battigalli and Siniscalchi (1999) proved the following result.

**Lemma 1.** Set  $\Delta^{\mathcal{F}}(\Omega)$  is compact metrizable.

Next, consider two compact metrizable spaces  $\Omega_1, \Omega_2$ , their product  $\Omega = \Omega_1 \times \Omega_2$  and a collection  $\mathcal{F} \subseteq \mathcal{B}(\Omega_1)$  of clopen subsets of  $\Omega_1$ . The “cylinders”

$$\mathcal{F}' = \{F \times \Omega_2 \subseteq \Omega : F \in \mathcal{F}\} \subseteq \mathcal{B}(\Omega)$$

form a family of clopen subsets of  $\Omega$ ; therefore,  $(\Omega, \mathcal{F}')$  is a conditional measurable space. We can marginalize CPSs defined over product spaces through the map  $\text{marg}_{\Omega_1} : \Delta^{\mathcal{F}'}(\Omega) \rightarrow \Delta^{\mathcal{F}}(\Omega_1)$  defined by

$$\text{marg}_{\Omega_1}(\mu)(E|F) = \mu(E \times \Omega_2 | F \times \Omega_2)$$

for all  $E \in \mathcal{B}(\Omega_1)$  and  $F \in \mathcal{F}$ , which is clearly continuous.

### 3. Multi-stage game form

A multi-stage game form is a mathematical object that encodes the rules of interaction in a game: the set of players (roles), and, for each player, his information and feasible actions at each stage as well as the outcomes of each feasible path of play. Here we consider a simplified version with *observable actions*. The game proceeds through stages. In each stage, the set of feasible actions of each player may depend on the history of past actions, which is public information. A player is inactive when his set of feasible actions is a singleton. In each stage, all active players move simultaneously.

A **multi-stage game tree** is a mathematical structure

$$(I, (A_i)_{i \in I}, \bar{H})$$

comprising the following elements:

- $I$  is the *finite* set of players;
- $A_i$  is a *finite* set of **actions** of player  $i \in I$ , and we let  $A = \prod_{i \in I} A_i$  denote the set of action profiles;
- $\bar{H} \subseteq \bigcup_{n \in \mathbb{N}_0} A^n$  denotes the *finite* set of **feasible histories**  $h = ((a_{i,s})_{i \in I})_{s=1}^n$ . In particular,  $\bar{H}$  has a *tree* structure: every prefix of a sequence in  $\bar{H}$  (including the empty sequence  $\emptyset$ ) belongs to  $\bar{H}$  as well. Thus, histories in  $\bar{H}$  correspond to nodes of the game tree and  $\emptyset$  is the root. Set  $\bar{H}$  is partitioned into the set of non-terminal histories  $H$  and terminal histories  $Z$ . Each  $z \in Z$  is a complete description of the *actual* behavior of players from the beginning to the end of the game.

For every  $h \in H$ , let

$$A(h) = \{a \in A : (h, a) \in \bar{H}\},$$

denote the set of feasible action profiles after  $h$ .<sup>14</sup> For every  $i \in I$ ,  $A_i(h)$  is the projection of  $A(h)$  onto  $A_i$  and  $A_{-i}(h)$  is similarly defined. The set of **personal histories** of player  $i$  is

<sup>14</sup> Note that, for each  $z \in Z$ ,  $A(z) = \emptyset$ .

$$H_i = \bar{H} \cup \{(h, a_i) \in H \times A_i : a_i \in A_i(h)\}.$$

In words,  $(h, a_i)$  represents the interim information of player  $i$  as soon as he has chosen action  $a_i$  given  $h$  and before he obtains information about the actions simultaneously chosen by the coplayers. Such interim histories are important to model what players believe about the consequences of their choices.

The natural precedence relation on  $\bar{H}$  is denoted by  $\preceq$ , that is, for all  $h, h' \in \bar{H}$ , we write  $h \preceq h'$  if and only if  $h$  is a prefix of  $h'$  (the irreflexive part of  $\preceq$  is denoted by  $<$ ). For each  $i \in I$ ,  $\preceq$  can be extended to a corresponding precedence relation  $\leq$  on  $H_i$  in an obvious way.<sup>15</sup> For every player  $i \in I$  and personal history  $h_i \in H_i$ ,

$$Z(h_i) = \{z \in Z : h_i \leq z\}$$

denotes the set of paths consistent with  $h_i$ . For each  $h \in H$ ,  $a_{-i} \in A_{-i}(h)$ , the subsets  $Z(h, a_i)$  and  $Z(h, a_{-i})$  of  $Z$  are similarly defined.

To complete the description of the game form we need to specify the relation between feasible actions and outcomes. In particular, we assume that the former and the latter may be mediated by players' personal traits such as intelligence, strength, or in general any kind of skills, according to natural laws. For instance, the exact consequences implied by a player doing a physical task may depend on his level of strength.

Formally, we parametrize the personal traits of each player  $i$  through a compact metrizable space  $\Theta_i$ . We assume that each player  $i$  knows his personal traits  $\theta_i \in \Theta_i$ , but is uncertain about the personal traits of coplayers,  $\theta_{-i} \in \Theta_{-i} = \prod_{j \in I \setminus \{i\}} \Theta_j$ . Players' personal traits  $\theta \in \Theta = \prod_{i \in I} \Theta_i$  are utility-relevant insofar they may influence both the outcomes and the preferences of players over outcomes and beliefs.<sup>16</sup> Thus, we assume that every player forms (potentially utility-relevant) beliefs about his coplayers' personal traits.<sup>17</sup>

We consider a metric space of outcomes  $Y$  and a continuous outcome function  $\pi : Z \times \Theta \rightarrow Y$  that maps each path of play and profile of personal traits to the corresponding outcome.<sup>18</sup> With this, we can define a **multi-stage game form** as a mathematical structure

$$\Gamma = \langle I, (A_i, \Theta_i)_{i \in I}, \bar{H}, \pi \rangle$$

As explained above, even though the “written” rules of the game do not explicitly refer to personal traits, natural laws may imply that outcomes depend on personal traits given players' actions.

We illustrate our framework and results with repeated reference to the game form depicted in Fig. 2.

**Example 1.** Ann ( $a$ ), Bob ( $b$ ), and Chloe ( $c$ ) play according to the following rules: First Ann and Bob choose independently and simultaneously between, respectively, Up or Down ( $U$  or  $D$ ), and Left or Right ( $L$  or  $R$ ). At this first stage, Chloe can only Wait ( $W$ ). (Waiting actions are not shown explicitly in the graphical representation.) If  $(D, L)$  is selected, Chloe can choose either a Non-selfish or a Selfish action ( $N$  or  $S$ ), while Ann and Bob can only Wait for the end of the interaction. Every other action pair of Ann and Bob terminates the interaction. The personal features of Ann and Bob are common knowledge, while the personal features of Chloe are summarized by a parameter  $\theta_c \in [\underline{\theta}_c, \bar{\theta}_c] \subseteq \mathbb{R}_+$  that measures Chloe's sensitivity to a given emotion (e.g., guilt, or anger) known only to Chloe. Outcomes are profiles of monetary payoffs, that is,  $Y = \mathbb{R}^I$ , where  $I = \{a, b, c\}$  is the player set, and  $\pi = (\pi_i)_{i \in I} : Z \times \Theta \rightarrow \mathbb{R}^I$ , where  $z \mapsto \pi_i(z, \theta)$  denotes the monetary payoff of  $i$  at terminal history  $z$  and is independent of  $\theta$ . Thus, in this case<sup>19</sup> the outcome depends only on what players do, not on their personal features.<sup>20</sup> We consider two possible monetary payoffs for Chloe after  $S$ ,  $x \in \{-1, 1\}$ , and we leave  $\pi_c((D, R, W))$  unspecified (denoted  $*$ ) to emphasize that it does not play any role in the numerical examples based on this game form. Consistently with the intended interpretation of our framework, we assume that this game form (including the value of  $x$ ) is common knowledge. To sum up, we have:<sup>21</sup>

- ▶  $A_a = \{U, D, W\}$ ,  $A_b = \{L, R, W\}$ ,  $A_c = \{N, S, W\}$ ,
- ▶  $H = \{\emptyset, ((D, L, W))\}$ ,
- ▶  $Z = \{((U, L, W)), ((U, R, W)), ((D, R, W)), ((D, L, W)), (W, W, S)\}$ ,
- ▶  $\Theta$  is isomorphic to  $[\underline{\theta}_c, \bar{\theta}_c]$ ,

and  $\pi$  is specified in the picture. One can derive from these primitive elements the sets of feasible actions and of personal histories of each player. For example,  $A_a(\emptyset) = \{U, D\}$ ,  $A_a((D, L, W)) = \{W\}$ , and  $H_a = H \cup Z \cup \{(U), (D)\}$ , where  $(U)$  and  $(D)$  are the interim histories associated with Ann's feasible actions at the root. ▲

<sup>15</sup> In particular, for  $h \in H$ ,  $a_i \in A_i(h)$ , and  $h' \in \bar{H}$ ,  $(h, a_i) \leq h'$  if and only if  $h < h'$  and  $(h, (a_i, a_{-i})) \leq h'$  for some  $a_{-i} \in A_{-i}(h)$ .

<sup>16</sup> See Section 5.

<sup>17</sup> See Section 4.

<sup>18</sup> To ease notation, we do not distinguish between traits affecting outcomes from traits that only affect preferences. Therefore, in some cases function  $\pi$  may be independent of  $\theta$ , as in the heuristic example of the Introduction.

<sup>19</sup> As in most experimental games.

<sup>20</sup> Payoff profiles are shown in alphabetical order from top to bottom.

<sup>21</sup> For the sake of clarity, we allow for redundant parentheses, as we have those of action profiles, those of sequences of action profiles (histories), and those of functions. To ease notation, in the following examples we omit redundant parentheses and the “waiting” actions.

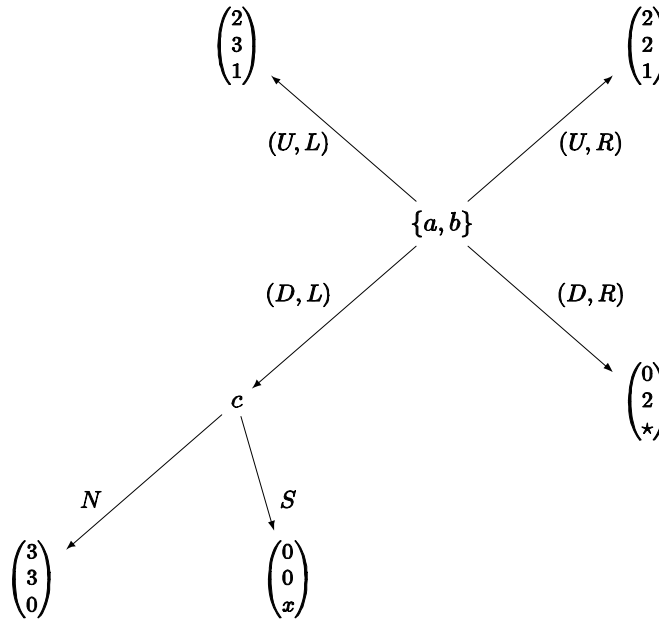


Fig. 2. A 3-person game form.

4. Beliefs

In this section we first analyze systems of conditional beliefs about paths and personal features of coplayers, focusing on a key independence property. Next we analyze the hierarchies of conditional beliefs.

4.1. Conditional beliefs, subjective plans and Own-action independence of beliefs

For a fixed player  $i \in I$ , we posit an abstract compact metrizable space  $T_{-i}$  that we interpret as the set of profiles of possible personal features of the coplayers. Specifically, each  $t_{-i} \in T_{-i}$  is interpreted as a description of the personal traits of the coplayers and of how they think, that is, what they think at the beginning of the game and how they update their beliefs upon receiving information about past play. For example, we may have  $T_{-i} = \Theta_{-i}$ . Later we will provide a constructive definition of this abstract space, which right now is not necessary.

The (abstract) uncertainty space of player  $i$  is the product space  $\Omega_{-i} = Z \times T_{-i}$ , and the corresponding collection of conditioning events is

$$\mathcal{F} = \{Z(h_i) \times T_{-i}\}_{h_i \in H_i} \subseteq \mathcal{B}(\Omega_{-i}).$$

Since  $Z$  is finite, each  $\Omega_{-i}(h_i) = Z(h_i) \times T_{-i}$  is clopen. Given that  $\mathcal{F}$  is isomorphic to  $H_i$ , let  $(\Omega_{-i}, H_i)$  denote the corresponding conditional measurable space. The space of CPSs on  $(\Omega_{-i}, H_i)$  is denoted by  $\Delta^{H_i}(\Omega_{-i})$ . For simplicity, we write conditional beliefs as  $\mu_i(\cdot|h_i)$ . Furthermore, for any CPS  $\mu_i$  on  $(\Omega_{-i}, H_i)$ , we introduce a simplified notation for marginal conditional probabilities summarized by the following Table 1, where  $h \in H$ ,  $h_i, h'_i \in H_i$  such that  $h_i \preceq h'_i$ ,  $E_{-i} \subseteq T_{-i}$  is measurable, and  $(a_i, a_{-i}) \in A(h)$ :

Table 1  
Marginal conditional probabilities.

Notation	Definition
$\mu_i(h'_i h_i)$	$\mu_i(Z(h'_i) \times T_{-i} h_i)$
$\mu_i(E_{-i} h_i)$	$\mu_i(Z \times E_{-i} h_i)$
$\mu_i(a_i, a_{-i} h)$	$\mu_i(Z(h, (a_i, a_{-i})) \times T_{-i} h)$
$\mu_i(a_i h)$	$\mu_i(Z(h, a_i) \times T_{-i} h) = \sum_{a'_{-i} \in A_{-i}(h)} \mu_i(a_i, a'_{-i} h)$
$\mu_i(a_{-i} h)$	$\mu_i(Z(h, a_{-i}) \times T_{-i} h) = \sum_{a'_i \in A_i(h)} \mu_i(a'_i, a_{-i} h)$

Define

$$\hat{\sigma}_i(\mu_i) = ((\mu_i(a_i|h))_{a_i \in A_i(h)})_{h \in H} \in \prod_{h \in H} \Delta(A_i(h))$$

and note that  $\hat{\sigma}_i(\mu_i)$  corresponds to a **behavior strategy** of  $i$ . We interpret these conditional probabilities as the subjective **plan** of  $i$  conditional on each possible history. Similarly,

$$\hat{\sigma}_{-i}(\mu_i) = \left( (\mu_i(a_{-i}|h))_{a_{-i} \in A_{-i}(h)} \right)_{h \in H} \in \prod_{h \in H} \Delta(A_{-i}(h))$$

corresponds to a **correlated behavior strategy** for player  $i$ 's coplayers. We interpret these conditional probabilities as the subjective belief of  $i$  about the behavior of his coplayers conditional on each possible history. Let  $\Sigma_i$  and  $\Sigma_{-i}$  respectively denote the set of subjective plans of  $i$  and the set of his beliefs about the behavior of his coplayers. Let  $\sigma_i$  (resp.  $\sigma_{-i}$ ) denote a generic element of  $\Sigma_i$  (resp.  $\Sigma_{-i}$ ) and let  $\hat{\sigma}_i(\mu_i)$  (resp.  $\hat{\sigma}_{-i}(\mu_i) \in \Sigma_{-i}$ ) denote the element of  $\Sigma_i$  (resp.  $\Sigma_{-i}$ ) derived from  $\mu_i$ . As we shall see in the following sections, the deterministic plans of  $i$  (i.e., those plans that assign probability 1 to a single action at each history) have a particular relevance in our analysis. Let  $S_i \subseteq \Sigma_i$  denote the set of deterministic plans of  $i$ , with generic element  $s_i$ .<sup>22</sup> In particular, with a slight abuse of notation, for each  $h \in H$ , let  $s_i(h)$  denote the unique action prescribed by a deterministic plan  $s_i$  at  $h$ . Note that each plan  $\sigma_i \in \Sigma_i$  can also be seen as a mixture over deterministic plans. Indeed, it will be useful to define the “support” of a plan  $\sigma_i \in \Sigma_i$  as

$$\text{supp}\sigma_i = \{s_i \in S_i : \forall h \in H, \sigma_i(s_i(h)|h) > 0\}.$$

As a final piece of notation, for all  $h_i, h'_i \in H_i$  with  $h_i \preceq h'_i$ , let  $\mathbb{P}_{\sigma_i, \sigma_{-i}}(h'_i|h_i)$  denote the conditional probability of  $h'_i$  given  $h_i$  derived from the behavior strategies  $(\sigma_i, \sigma_{-i}) \in \Sigma_i \times \Sigma_{-i}$  in the usual way.

The definition of conditional belief as expressed in Definition 2 may represent the system of beliefs of an external observer that obtains the same information as player  $i$ . But, arguably, reasonable beliefs of player  $i$  should satisfy a further condition: *what  $i$  believes about his coplayers' features and simultaneous actions is independent of his own actions*. Indeed, there is no objective causality that links the behavior of player  $i$  at a certain history with the simultaneous behavior of his coplayers or their personal features. We next formalize this property and then illustrate it through our running example.

**Definition 3.** A CPS  $\mu_i$  on  $(\Omega_{-i}, H_i)$  satisfies *own-action independence* (OAI) if

$$\mu_i(Z(h, (a_i, a_{-i})) \times E_{-i}|h, a_i) = \mu_i(Z(h, (a'_i, a_{-i})) \times E_{-i}|h, a'_i)$$

for all  $h \in H$ ,  $a_i, a'_i \in A_i(h)$ ,  $a_{-i} \in A_{-i}(h)$  and measurable  $E_{-i} \subseteq T_{-i}$ .<sup>23</sup>

The set of CPSs of  $i$  satisfying OAI is denoted by  $\Delta_i^{H_i}(\Omega_{-i})$ .

**Example 2.** Consider the game form depicted in Fig. 2. Let us focus on Ann's initial beliefs on  $Z \times \Theta_c$ , and those conditional on her action  $D$  and action pair  $(D, L)$ . Suppose that at the root Ann plans to go Up, she expects Bob to go Left with probability 0.5, and she assigns probability 0.5 to each one of the two extreme types  $\underline{\theta}_c$  and  $\bar{\theta}_c$  of Chloe. Thus,  $\mu_a((U, L, \bar{\theta}_c)|\emptyset) = 0.25$  and  $\mu_a((D, a_b, \theta_c)|\emptyset) = 0$  for every  $(a_b, \theta_c) \in A_b(\emptyset) \times \Theta_c$ . Condition OAI implies  $\mu_a(L, \bar{\theta}_c|\emptyset) = \mu_a(L, \bar{\theta}_c|D)$ , because what Ann believes about the simultaneous action of Bob and the type of Chloe is independent of what she does, even when—as in this case—Ann “surprises herself” by taking an unplanned action. Yet, OAI does not require independence across coplayers: in particular, it is possible that  $\mu_a(L, \bar{\theta}_c|\emptyset) \neq 0.5 \times 0.5$  and  $\mu_a(\bar{\theta}_c|(D, L)) \neq 0.5$ , i.e., that in Ann's eyes there is (spurious) correlation between Bob's behavior and Chloe's type.  $\blacktriangle$

The following proposition characterizes OAI.

**Proposition 1.** Consider a CPS  $\mu_i$  on  $(\Omega_{-i}, H_i)$ . The following are equivalent:

- i)  $\mu_i$  satisfies OAI;
- ii) for all  $h \in H$ ,  $(a_i, a_{-i}) \in A(h)$ , and measurable  $E_{-i} \subseteq T_{-i}$ ,

$$\mu_i(Z(h, a_{-i}) \times E_{-i}|h) = \mu_i(Z(h, (a_i, a_{-i})) \times E_{-i}|h, a_i).$$

OAI implies that a CPS of player  $i$  is made of two independent parts,  $i$ 's beliefs about his own behavior and his beliefs about the coplayers' behavior and personal features. In the Appendix, we also show that OAI is equivalent to a factorization of  $\mu_i \in \Delta_i^{H_i}(\Omega_{-i})$  in a behavior strategy  $\hat{\sigma}_i(\mu_i) \in \Sigma_i$  for  $i$  and a pair

<sup>22</sup> Note that, mathematically, the set of deterministic plans just defined coincides with the standard set of pure strategies for an extensive-form game. Indeed, we are denoting them with the usual notation  $s_i \in S_i$ . However, according to our interpretation, each  $s_i \in S_i \subseteq \Sigma_i$  is a subjective plan which represents player  $i$ 's intentions rather than actual behavior. In fact, the latter is described only by play paths  $z \in Z$ .

<sup>23</sup> The previous condition holds vacuously for terminal information sets.



$$(\zeta_{-i}, \eta_i) = \left( \left( \zeta_{-i}^{t_{-i}} \right)_{t_{-i} \in T_{-i}}, (\eta_{i,h})_{h \in \bar{H}} \right) \in \Sigma_{-i}^{T_{-i}} \times \Delta(T_{-i})^{\bar{H}}$$

of features-dependent behavior strategies of coplayers and a system of marginal beliefs of  $i$  over  $T_{-i}$ . In words, each  $\zeta_{-i}^{t_{-i}}$  describes the belief of  $i$  about the behavior of his coplayers when their personal features are  $t_{-i}$ , and each  $\eta_{i,h}$  describes the beliefs of  $i$  about his coplayers' personal features at every history. In particular, the correlated behavior strategy  $\hat{\sigma}_{-i}(\mu_i) \in \Sigma_{-i}$  defined above is obtained as

$$\hat{\sigma}_{-i}(\mu_i)(a_{-i}|h) = \int_{T_{-i}} \zeta_{-i}^{t_{-i}}(a_{-i}|h) \eta_{i,h}(dt_{-i})$$

for all  $h \in H$  and  $a_{-i} \in A_{-i}(h)$ . Clearly, in two-person game forms with observable actions, each  $\hat{\sigma}_{-i}(\mu_i)$  is a proper behavior strategy for the (unique) coplayer of  $i$ .

Relying on Lemma 1, we can prove the following result.

**Lemma 2.** Set  $\Delta_i^{H_i}(\Omega_{-i})$  is compact metrizable.

To ease the exposition, we refer to CPS that satisfy the OAI property with the acronym ICPS.

**Definition 4.** We say that an ICPS  $\mu_i \in \Delta_i^{H_i}(\Omega_{-i})$  strongly believes event  $E_{-i} \in \mathcal{B}(\Omega_{-i})$  if, for every  $h \in H$ ,

$$\Omega_{-i}(h) \cap E_{-i} \neq \emptyset \Rightarrow \mu_i(E_{-i}|h) = 1.$$

We say that  $\mu_i \in \Delta_i^{H_i}(\Omega_{-i})$  strongly believes  $(E_{-i}^1, \dots, E_{-i}^n) \in \mathcal{B}(\Omega_{-i})^n$  if  $\mu_i$  strongly believes  $E_{-i}^k$  for every  $k \in \{1, \dots, n\}$ .<sup>24</sup>

The following result adapts Lemma 3 of Battigalli and Tebaldi (2019).<sup>25</sup> It is crucial in the proof of the main theorem of this article.

**Lemma 3.** Let  $\Omega_{-i} = Z \times T_{-i}$  as above and let  $X_{-i}$  be compact metrizable. Fix a decreasing chain  $(E_{-i}^1, \dots, E_{-i}^n)$  in  $\Omega_{-i} \times X_{-i}$  and let  $\text{proj}_{\Omega_{-i}} E_{-i}^m$  be measurable for each  $m \in \{1, \dots, n\}$ . For each  $\mu_i \in \Delta_i^{H_i}(\Omega_{-i})$  that strongly believes  $(\text{proj}_{\Omega_{-i}} E_{-i}^1, \dots, \text{proj}_{\Omega_{-i}} E_{-i}^n)$  there exists  $\nu_i \in \Delta_i^{H_i}(\Omega_{-i} \times X_{-i})$  that strongly believes  $(E_{-i}^1, \dots, E_{-i}^n)$  and satisfies  $\text{marg}_{\Omega_{-i}} \nu_i = \mu_i$ .

#### 4.2. Hierarchies of conditional beliefs

For each player  $i \in I$ , the space of **primitive uncertainty** of  $i$  is  $\Omega_{-i}^0 = Z \times \Theta_{-i}$ , that is, player  $i$  is first of all uncertain about how the game is going to be played ( $z \in Z$ ) and of the personal traits of the coplayers ( $\theta_{-i} \in \Theta_{-i}$ ). The set  $\Omega_{-i}^0$  is compact metrizable because  $Z$  is finite and  $\Theta_{-i}$  is compact metrizable. It is immediate to see that  $\Omega_{-i}^0$  is a particular case of the (abstract) uncertainty space  $\Omega_{-i}$  introduced in Section 4.1: it is enough to let  $T_{-i} = \Theta_{-i}$ . With this, the conditional measurable space  $(\Omega_{-i}^0, H_i)$  is well defined and we can consider ICPSs in  $\Delta_i^{H_i}(\Omega_{-i}^0)$ .

For all  $i \in I$ , hierarchies of ICPSs are recursively defined as follows:

- $\Omega_{-i}^0 = Z \times \Theta_{-i}$ ;  $M_{i,1} = M_i^1 = \Delta_i^{H_i}(\Omega_{-i}^0)$ ,
- $\Omega_{-i}^k = \Omega_{-i}^{k-1} \times \prod_{j \in I \setminus \{i\}} M_{j,k}$ ;  $M_{i,k+1} = \Delta_i^{H_i}(\Omega_{-i}^k)$ ;  $M_i^{k+1} = \prod_{m=1}^{k+1} M_{i,m}$  ( $k \in \mathbb{N}$ );

where, for all  $k \in \mathbb{N}$ ,  $\Omega_{-i}^k$  is the  **$k$ -th order uncertainty space** of  $i$ ,  $M_{i,k}$  is the space of  **$k$ -th order ICPSs** of  $i$  and  $M_i^k$  is the space of  **$k$ -th order hierarchies** of ICPSs of  $i$ . A generic element  $\mu_i^k = (\mu_{i,m})_{m=1}^k$  in  $M_i^k$  is a sequence of length  $k$  of ICPSs defined on uncertainty spaces of increasing orders. Moreover, for all  $k \in \mathbb{N}$ , it can be checked that

$$\Omega_{-i}^k = Z \times \Theta_{-i} \times \prod_{m=1}^k \prod_{j \in I \setminus \{i\}} M_{j,m}.$$

For each  $i \in I$  and  $k \in \mathbb{N}$ , let

<sup>24</sup> We can show that OAI implies that if the condition displayed above holds for all  $h \in \bar{H}$ , then it also holds for all  $h_i \in H_i$ .

<sup>25</sup> In particular, we rely on a generalization of their result proved in the working paper version of the article (IGIER w.p. 609).

$$M_{-i}^k = \prod_{m=1}^k \prod_{j \in I \setminus \{i\}} M_{j,m} \text{ and } T_{-i}^k = \Theta_{-i} \times M_{-i}^k$$

respectively denote the spaces of  $k$ -th order hierarchies and personal features of coplayers; with this, each  $k$ -th order uncertainty space  $\Omega_{-i}^k = Z \times T_{-i}^k$  has the same structure of the abstract space  $\Omega_{-i}$  introduced in Section 4.1. By repeated applications of Lemma 2 and Tychonoff's Theorem, one can show that each  $M_i^k$  is compact metrizable. Thus, for every  $k \in \mathbb{N}$ ,  $T_{-i}^k$  is compact metrizable and  $(\Omega_{-i}^k, H_i)$  is a well defined conditional measurable space.<sup>26</sup>

An **infinite hierarchy of ICPSs** of player  $i$  is a denumerable sequence  $\mu_i^\infty = (\mu_{i,k})_{k \in \mathbb{N}} \in M_i^\infty$ , where  $M_i^\infty = \prod_{k \in \mathbb{N}} M_{i,k}$ . For every  $i \in I$ ,  $M_i^\infty$  is compact metrizable, and so are  $M_{-i}^\infty = \prod_{j \in I \setminus \{i\}} M_j^\infty$  and  $M^\infty = \prod_{i \in I} M_i^\infty$ . The **total state space** given the game form  $\Gamma$  is

$$\Omega^\infty = Z \times \Theta \times M^\infty.$$

Each state  $\omega^\infty = (z, \theta, \mu^\infty)$  in  $\Omega^\infty$  is a complete description of the actual behavior of players, of their personal traits, and of their systems of conditional beliefs of all orders. Since player  $i$  is assumed to know both his personal traits  $\theta_i$  and his infinite hierarchy  $\mu_i^\infty$ , the **total uncertainty space** of  $i$  is

$$\Omega_{-i}^\infty = Z \times \Theta_{-i} \times M_{-i}^\infty.$$

The **spaces of personal features** of player  $i$  and coplayers  $-i$  are, respectively,  $T_i^\infty = \Theta_i \times M_i^\infty$  and  $T_{-i}^\infty = \Theta_{-i} \times M_{-i}^\infty$ . With this, we can write the total uncertainty space of  $i$  as  $\Omega_{-i}^\infty = Z \times T_{-i}^\infty$ , and  $(\Omega_{-i}^\infty, H_i)$  is a well defined conditional measurable space.

### 4.3. Belief coherence

So far, we did not impose the requirement that beliefs of different orders in a hierarchy are mutually consistent, i.e., that they assign the same probabilities to events of lower order of uncertainty, such as events about behavior. Next we consider hierarchies satisfying this requirement, that we interpret as a cognitive rationality condition.<sup>27</sup>

**Definition 5.** Fix an infinite belief hierarchy  $\mu_i^\infty \in M_i^\infty$ . We say that  $\mu_i^\infty$  is **coherent** if, for all  $k \in \mathbb{N}$ , and  $h_i \in H_i$ ,

$$\text{marg}_{\Omega_{-i}^{k-1}} \mu_{i,k+1}(\cdot|h_i) = \mu_{i,k}(\cdot|h_i). \tag{2}$$

For each  $i \in I$ , we let  $C_i^\infty$  denote the subset of coherent hierarchies in  $M_i^\infty$  and let  $C^\infty = \prod_{i \in I} C_i^\infty$ . For each  $n \in \mathbb{N}$ ,  $C_i^n$  and  $C^n$  are similarly defined, i.e., condition (2) must hold for all  $k \leq n$ . The following technical result implies that  $C_i^\infty$  has the same topological properties of  $M_i^\infty$ .

**Lemma 4.** Set  $C_i^\infty$  is compact metrizable.

This easily follows from the fact that, for all  $k \in \mathbb{N}$  and  $h_i \in H_i$ , the map

$$\mu_{i,k+1}(\cdot|h_i) \mapsto \text{marg}_{\Omega_{-i}^{k-1}} \mu_{i,k+1}(\cdot|h_i)$$

is continuous (see for example Aliprantis and Border, 2006, Chapter 15).

The following result adapts Proposition 1 of Battigalli and Siniscalchi (1999) to the current setup.<sup>28</sup>

**Proposition 2.** There exists a **canonical homeomorphism**  $g_i : C_i^\infty \rightarrow \Delta_i^{H_i}(\Omega_{-i}^\infty)$  such that for all  $\mu_i^\infty \in C_i^\infty$  and  $k \in \mathbb{N}$ ,  $\text{marg}_{\Omega_{-i}^{k-1}} g_i(\mu_i^\infty) = \mu_{i,k}$ .

<sup>26</sup> The countable collection of clopen subsets of  $\Omega_{-i}^k$  is given by

$$\left\{ Z(h_i) \times T_{-i}^k \right\}_{h_i \in H_i} \subseteq \mathcal{B}(\Omega_{-i}^k),$$

which is isomorphic to  $H_i$  for all  $k \in \mathbb{N}$ .

<sup>27</sup> Note, we did already impose some cognitive rationality requirements: the chain rule and OAI. We did this only for brevity. We could allow for more deviations from cognitive rationality at arbitrary states of the world, requiring all the relevant consistency conditions only at states where players are rational.

<sup>28</sup> Proposition 1 of Battigalli and Siniscalchi (1999) considers an abstract conditional measurable space, not one based on a game; therefore, the game-based OAI condition used in the definition of  $C_i^\infty$  and  $\Delta_i^{H_i}(\Omega_{-i}^\infty)$  cannot be expressed. The game-theoretic applications of the second part of that article instead impose a requirement in the spirit of OAI as part of the definition of rationality.

Finally, it is straightforward to define the space of **coherent  $k$ -th order hierarchies** of ICPSs which is denoted by  $C_i^k$ .

**5. Psychological games and rationality**

Next we introduce belief-dependent utilities and append them to the game form to obtain a psychological game. Fix a multi-stage game form  $\Gamma$ , a  **$k$ -th order psychological game** based on  $\Gamma$  is a structure

$$(\Gamma, v) = (\Gamma, (v_i)_{i \in I})$$

where, for every  $i \in I$ , the **psychological utility function**  $v_i : Y \times \Theta \times M^k \rightarrow \mathbb{R}$  is *continuous*.<sup>29</sup> Let  $PG_k$  denote the class of  $k$ -th order psychological games and let  $(\Gamma, v)$  denote a generic element of this class. Note that in many interesting cases  $v_i$  depends only on  $\theta_i \in \Theta_i$  (which parametrizes player  $i$ 's preferences over outcomes and beliefs), not on  $\theta_{-i}$ . However, we allow for dependence of each  $v_i$  on the entire profile  $\theta \in \Theta$ . With this, we are also able to consider models of interdependent preferences as in Gul and Pesendorfer (2016).<sup>30</sup>

For simplicity, we will always consider the implied utility function defined over terminal histories and personal features

$$u_i : Z \times \Theta \times M^k \rightarrow \mathbb{R},$$

$$(z, \theta, \mu^k) \mapsto v_i(\pi(z, \theta), \theta, \mu^k),$$

obtained as composition of the outcome function  $\pi$  with  $v_i$ . With this, the **utility-relevant state space** of each player is  $\Omega^k = Z \times \Theta \times M^k$ . Note that the total state space can be factorized as follows:

$$\Omega^\infty = \Omega^k \times \prod_{i \in I} \prod_{m \geq k} \Delta_i^{H_i}(\Omega_{-i}^m) = Z \times \prod_{i \in I} \left( T_i^k \times \prod_{m > k} M_i^m \right), \tag{3}$$

that is, each state  $\omega^\infty \in \Omega^\infty$  can be seen as a pair  $(\omega^k, \omega^{>k})$  such that only  $\omega^k \in \Omega^k$  is utility-relevant. Since player  $i$  knows his own personal features  $(\theta_i, \mu_i^k) \in \Theta_i \times M_i^k = T_i^k$ , his **utility-relevant uncertainty space** is  $\Omega_{-i}^k = Z \times T_{-i}^k$ .

A rational player  $i$  has to consult his  $(k + 1)$ -order ICPS  $\mu_{i,k+1}$  in order to take conditional expectations of his psychological utility. If  $i$  has observed  $h \in H$ , for each  $a_i \in A_i(h)$ , he computes the corresponding expected utility given  $\theta_i$  and  $\mu_{i,k+1}$ , that is,

$$\bar{u}_{i,h}(a_i, \theta_i, \mu_{i,k+1}) = \int_{\Omega_{-i}^k} u_i(\theta_i, \mu_i^k, \omega_{-i}^k) \mu_{i,k+1}(d\omega_{-i}^k | h, a_i). \tag{4}$$

Hence, for every  $i \in I$ , we have a vector

$$(\bar{u}_{i,h} : A_i(h) \times \Theta_i \times M_{i,k+1} \rightarrow \mathbb{R})_{h \in H}$$

of **psychological decision utilities** defined as in (4).

**Remark 1.** For every  $i \in I$  and  $h \in H$ ,  $\bar{u}_{i,h}$  is continuous.

In the rest of the paper, we only refer to such decision utilities to define players' rationality and derive our results. Note that we could have considered the vector of utilities  $(\bar{u}_{i,h})_{(i,h) \in I \times H}$  as the primitive element of the analysis, rather than deriving it from the psychological utility defined over terminal paths.<sup>31</sup>

Given  $h \in H$ , a rational player  $i$  with personal traits  $\theta_i$  and beliefs  $\mu_{i,k+1}$  solves the following problem:

$$\max_{a_i \in A_i(h)} \bar{u}_{i,h}(a_i, \theta_i, \mu_{i,k+1}).$$

With this, for every  $h \in H$ , we define the **local best-reply correspondence** of  $i$  at  $h$  as

$$r_{i,h} : \Theta_i \times M_{i,k+1} \rightrightarrows A_i(h),$$

$$(\theta_i, \mu_{i,k+1}) \mapsto \arg \max_{a_i \in A_i(h)} \bar{u}_{i,h}(a_i, \theta_i, \mu_{i,k+1}).$$

**Remark 2.** For every  $i \in I$  and  $h \in H$ ,  $r_{i,h}$  is non-empty valued and upper hemicontinuous.

<sup>29</sup> As in Battigalli and Dufwenberg (2009), we take continuity of psychological utilities as a maintained assumption. In their analysis of static psychological games, Jagau and Perea (2017) relax the continuity assumption.

<sup>30</sup> See also Levine (1998).

<sup>31</sup> Such reduced form approach is more general because it does not require that decision utilities are derived from an overall utility over terminal nodes. As explained in Battigalli et al. (2019a) this allows to model some belief-dependent action tendencies like being aggressive when frustrated.

In the following section, we are going to define rational planning of  $i$  given the personal features  $(\theta_i, \mu_{i,k+1})$  as a consistency condition between the local best replies  $(r_{i,h}(\theta_i, \mu_{i,k+1}))_{h \in H}$  and the plan (behavior strategy)  $\hat{\sigma}_i(\mu_{i,1})$  of  $i$  derived by his first-order beliefs. Recall that in Section 4 we have defined the subjective plan of  $i$  implied by his beliefs as:

$$\hat{\sigma}_i(\mu_{i,1})(a_i|h) = \text{marg}_Z \mu_{i,1}(Z(h, a_i)|h)$$

for all  $h \in H$  and  $a_i \in A_i(h)$ .<sup>32</sup> Note that here we take the space of first-order beliefs as domain of the map  $\hat{\sigma}_i$ . This choice—although natural—is somewhat arbitrary, because one can derive a plan for  $i$  from a belief of any order. Yet, we will focus on plans of rational players, whose belief hierarchy is coherent and hence yields the same plan starting from beliefs of any order.

### 5.1. Rational planning and material consistency

Recall that player  $i$  has a coherent belief at  $\omega_i^\infty = (z, \theta_i, \mu_i^\infty) \in \Omega_i^\infty$  whenever  $\mu_i^\infty \in C_i^\infty$ , that is, whenever the prevailing state of the world  $\omega_i^\infty$  belongs to

$$\Omega_i^{\infty,*} = \{(\bar{z}, \bar{\theta}_i, \bar{\mu}_i^\infty) \in \Omega_i^\infty : \bar{\mu}_i^\infty \in C_i^\infty\}$$

which is measurable. Event  $\Omega_i^{\infty,*}$  represents the statement “Player  $i$  has coherent beliefs.”

**Remark 3.** The set  $\Omega_i^{\infty,*}$  is compact.<sup>33</sup>

Our notion of *rationality* is given by the conjunction of several consistency conditions. We begin with rational planning which requires the subjective plan  $\hat{\sigma}_i(\mu_{i,1})$  of  $i$  to be immune to one-shot deviations given his personal features  $(\theta_i, \mu_{i,k+1})$ .

**Definition 6.** Player  $i$  is a **rational planner (RP)** at  $(z, \theta_i, \mu_i^\infty) \in \Omega_i^\infty$  if  $\mu_i^\infty \in C_i^\infty$  and

$$\hat{\sigma}_i(\mu_{i,1})(r_{i,h}(\theta_i, \mu_{i,k+1})|h) = 1 \tag{5}$$

for every  $h \in H$ .

In words, rational planning corresponds to a notion of *intra-personal equilibrium* among the selves (indexed by all the feasible histories  $h \in H$ ) of player  $i$ . Indeed, it requires that all the actions  $a_i$  in the support of the plan  $\hat{\sigma}_i(\mu_{i,1})(\cdot|h)$  maximize the local utility  $\bar{u}_{i,h}(\cdot, \theta_i, \mu_{i,k+1})$  which on turn depends on the first-order belief  $\mu_{i,1}$ , hence on the plan  $\hat{\sigma}_i(\mu_{i,1})$  itself.<sup>34</sup> Therefore,  $i$  is a rational planner if and only if his belief  $\mu_{i,k+1}$  satisfies the fixed-point condition in Eq. (5) for every  $h \in H$ . In Example 3 we illustrate this point by considering psychological preferences directly expressed through local decision utilities.

The event that player  $i$  is a rational planner is

$$RP_i = \{(z, \theta_i, \mu_i^\infty) \in \Omega_i^{\infty,*} : \forall h \in H, \hat{\sigma}_i(\mu_{i,1})(r_{i,h}(\theta_i, \mu_{i,k+1})|h) = 1\}.$$

**Lemma 5.** The set  $RP_i$  is non-empty and compact.

Another aspect of rationality is the consistency between planned behavior and actual behavior:

**Definition 7.** Player  $i$  is **materially consistent (MC)** at  $(z, \theta_i, \mu_i^\infty) \in \Omega_i^\infty$  if  $\mu_i^\infty \in C_i^\infty$  and, for all  $h \in H$ ,

$$h < z \implies \hat{\sigma}_i(\mu_{i,1})(a_{i,h}(z)|h) > 0,$$

where  $a_{i,h}(z) \in A_i(h)$  is the unique feasible action of  $i$  at  $h$  implied by  $z$ .

The corresponding event is

$$MC_i = \{(z, \theta_i, \mu_i^\infty) \in \Omega_i^{\infty,*} : \forall h \in H, h < z \implies \hat{\sigma}_i(\mu_{i,1})(a_{i,h}(z)|h) > 0\}.$$

Note that we can similarly define the event that player  $i$  is **strictly materially consistent** as

$$MC_i^* = \{(z, \theta_i, \mu_i^\infty) \in MC_i : \forall h \in H, \exists a_i \in A_i(h), \hat{\sigma}_i(\mu_{i,1})(a_i|h) = 1\}$$

where we require that player  $i$ 's subjective plan assigns probability 1 to his actual behavior.

<sup>32</sup> See Section 4 and Proposition 6 in the Appendix for more details.

<sup>33</sup> To see this, just notice that  $\Omega_i^{\infty,*} = Z \times \Theta_i \times C_i^\infty$ , where  $C_i^\infty$  is compact (Lemma 4).

<sup>34</sup> Note that here the coherence property between  $\mu_{i,1}$  and  $\mu_{i,k+1}$  is crucial.

**Lemma 6.** Set  $MC_i$  is nonempty and measurable,  $MC_i^*$  is nonempty and compact.

Rationality is therefore defined as follows.

**Definition 8.** Player  $i$  is **rational** at  $(z, \theta_i, \mu_i^\infty) \in \Omega_i^\infty$  if  $i$  is a rational planner and is materially consistent at  $(z, \theta_i, \mu_i^\infty)$ .

The event that player  $i$  is rational is denoted by  $R_i = RP_i \cap MC_i$ . Finally, we define the events “Every player is rational” and “Every coplayer of player  $i$  is rational,” respectively in  $\Omega^\infty$  and  $\Omega_{-i}^\infty$ , as:

$$R = \bigcap_{i \in I} (R_i \times T_{-i}^\infty) \text{ and } R_{-i} = \bigcap_{j \in I \setminus \{i\}} \left( R_j \times \prod_{t \in I \setminus \{i, j\}} T_t^\infty \right),$$

with the convention that, whenever we analyze 2-player games, the set  $\prod_{t \in I \setminus \{i, j\}} T_t^\infty$  is a singleton and  $R_{-i}$  is equal to  $R_j$ . Moreover, we define the sets  $R_i^*$ ,  $R^*$  and  $R_{-i}^*$  by replacing  $MC_i$  with  $MC_i^*$  in all the previous definitions.

**Remark 4.** Sets  $R_i$ ,  $R_{-i}$  and  $R$  are non-empty and measurable; sets  $R_i^*$ ,  $R_{-i}^*$  and  $R^*$  are compact.

Note that  $R_i^*$  is the event that  $i$  is rational and has a deterministic plan. This set can be empty. Indeed, there are cases in which only non-deterministic plans are consistent with rational planning.

We illustrate our notion of subjective rationality in a case where psychological preferences are not dynamically consistent.

**Example 3.** Consider the game form with monetary payoffs of Fig. 2 with  $x = -1$ , so that  $S$  is a costly punishment that Chloe can inflict on Ann and Bob. Suppose that Chloe, given  $(D, L)$ , may be affected by frustration and simple anger as in the theory of Battigalli et al. (2019b). In this particular case, her frustration is measured by her disappointment,<sup>35</sup> which induces a propensity to harm coplayers proportional to frustration. Specifically, Chloe’s decision utility function is

$$\bar{u}_{c,(D,L)}(a_c, \theta_c, \mu_{c,1}) = \pi_c((D, L), a_c) - \theta_c [\mathbb{E}_{\mu_{c,1}}(\pi_c | \emptyset)]^+ \sum_{j \in \{a,b\}} \pi_j((D, L), a_c),$$

where  $[\cdot]^+ = \max\{\cdot, 0\}$ . Suppose Chloe initially expects Bob to go Left with probability 1 and Ann to go Down with probability 0.5, so that  $\hat{\sigma}_{-c}(\mu_{c,1})((D, L) | \emptyset) = 0.5$ , and let  $p = \hat{\sigma}_c(\mu_{c,1})(S | (D, L))$  denote the probability with which Chloe plans to “punish” her coplayers. Then

$$\mathbb{E}_{\mu_{c,1}}(\pi_c | \emptyset) = 0.5 + 0.5px = \frac{1-p}{2} \geq 0$$

and the decision utilities of  $N$  and  $S$  are

$$\bar{u}_{c,(D,L)}(N, \theta_c, \mu_{c,1}) = -3\theta_c(1-p), \bar{u}_{c,(D,L)}(S, \theta_c, \mu_{c,1}) = -1$$

respectively. Given her beliefs about coplayers, Chloe cannot have a deterministic rational plan to punish, because  $p = 1$  implies that she wants to choose  $N$  given  $(D, L)$ :

$$\bar{u}_{c,(D,L)}(N, \theta_c, \mu_{c,1}) = 0 > -1 = \bar{u}_{c,(D,L)}(S, \theta_c, \mu_{c,1}).$$

Also, she has a deterministic rational plan not to punish ( $p = 0$ ) if and only if

$$\bar{u}_{c,(D,L)}(N, \theta_c, \mu_{c,1}) = -3\theta_c \geq -1 = \bar{u}_{c,(D,L)}(S, \theta_c, \mu_{c,1}),$$

i.e.,  $\theta_c \leq 1/3$ . Thus, if  $\theta_c > 1/3$ , Chloe has no deterministic rational plan. In this case, the non-deterministic rational plan solves the following indifference condition

$$3\theta_c(1-p) = 1,$$

which yields  $p = 1 - 1/(3\theta_c) \in (0, 1)$ . With this, Chloe plans rationally and is materially consistent, hence is rational, at both personal states  $((D, L), a_c, \theta_c, \mu_{c,1})$  ( $a_c \in \{N, S\}$ ) if  $\theta_c > 1/3$  and  $p = 1 - 1/(3\theta_c)$ , because in both states she plans to choose with positive probability a local best reply, which is also her action on the actual path. ▲

<sup>35</sup> In the theory of Battigalli et al. (2019b), disappointment is only an upper bound on frustration. Note also that, as mentioned earlier, in the theory of anger decision utilities are not simply derived as the expectations of terminal psychological utility.

In the previous example, dynamic inconsistency of preferences implies the non-existence of a rational deterministic plan. Next we provide a simple condition on psychological utility such that the induced preferences over actions are dynamically consistent, because they are represented by a *state-dependent subjective expected utility*.

**Definition 9.** We say that  $(\Gamma, v)$  satisfies **own-belief independence** (OBI) if, for all  $i \in I$ ,  $y \in Y$ ,  $\theta \in \Theta$  and  $\mu^k, \hat{\mu}^k \in M^k$ ,

$$\mu_{-i}^k = \hat{\mu}_{-i}^k \implies v_i(y, \theta, \mu^k) = v_i(y, \theta, \hat{\mu}^k).$$

Let  $PG_k^{OBI}$  denote the class of  $k$ -th order psychological games satisfying OBI. In words, a psychological game satisfies OBI if and only if, for each player  $i$  with traits  $\theta_i$ , the utility of outcomes depends on the unknown state  $(\theta_{-i}, \mu_{-i}^k)$ , that is, it can be written as  $v_{i,\theta_i}(y, \theta_{-i}, \mu_{-i}^k)$ . In particular,  $i$ 's utility is independent of his own plan  $\hat{\sigma}_i(\mu_{i,1})$ .<sup>36</sup> For example, when each  $v_i$  is given by the guilt-aversion formula (1) of Section 1, we obtain a game in  $PG_1^{OBI}$ . Models of image concerns—whereby  $i$ 's utility depends on  $-i$ 's terminal beliefs about unobserved traits or actions of  $i$ —provide another example of psychological preferences satisfying OBI.<sup>37</sup> Under OBI, we can rely on standard results about subjective expected utility maximization, such as the one-shot deviation principle and the existence of deterministic rational plans. In order to formally express such results we need a further piece of notation. For all  $h \in H$ ,  $\mu_{i,k+1} \in M_{i,k+1}$ , and  $\sigma_i \in \Sigma_i$  define  $\mu_{i,k+1}^{\sigma_i}$  as the  $k+1$  order ICPS obtained from  $\mu_{i,k+1}$  by substituting the original plan  $\hat{\sigma}_i(\mu_{i,1})$  implied by  $\mu_{i,k+1}$  with the arbitrary plan  $\sigma_i$ . Note that OAI implies that this factorization is well defined.<sup>38</sup> Also, for all  $h \in H$ ,  $\sigma_i \in \Sigma_i$ ,  $\theta_i \in \Theta_i$ , and  $\mu_{i,k+1} \in M_{i,k+1}$ , define

$$V_{i,h}(\sigma_i, \theta_i, \mu_{i,k+1}) = \int_{\Omega_{-i}^k} u_i(\theta_i, \mu_i^k, \omega_{-i}^k) \mu_{i,k+1}^{\sigma_i}(d\omega_{-i}^k | h).$$

**Remark 5.** If  $(\Gamma, v) \in PG_k^{OBI}$ , then player  $i$  is a rational planner at  $(z, \theta_i, \mu_i^\infty) \in \Omega_i^\infty$  if and only if  $\mu_i^\infty \in C_i^\infty$  and, for each  $h \in H$ ,

$$s_i \in \text{supp} \hat{\sigma}_i(\mu_{i,1}) \implies s_i^h \in \arg \max_{s'_i \in S_i(h)} V_{i,h}(s'_i, \theta_i, \mu_{i,k+1}),$$

where  $s_i^h$  is the deterministic plan that allows  $h$  and coincides with  $s_i$  at all histories that do not precede  $h$ , and  $S_i(h) \subseteq \Sigma_i$  is the set of deterministic plans that allow history  $h$ .

The previous remark is just a restatement adapted to the current framework of the standard equivalence between the one-shot deviation property and sequential optimality of strategies under subjective expected utility maximization.<sup>39</sup>

### 6. Strong belief and strong rationalizability

In this section, we adapt Battigalli and Siniscalchi's (2002) notion of strong belief to our framework and then provide the definition of strong rationalizability. We first focus on events within the utility-relevant space of uncertainty of player  $i$ . We say that  $(k+1)$  order belief  $\mu_{i,k+1} \in M_{i,k+1}$  **strongly believes** an event  $E_{-i} \subseteq \Omega_{-i}^k$  if

$$\Omega_{-i}^k(h) \cap E_{-i} \neq \emptyset \implies \mu_{i,k+1}(E_{-i} | h) = 1,$$

for all  $h \in H$ .

We are now ready to present the algorithm defining **strong rationalizability**, which is meant to capture the predictions for behavior and low-order beliefs implied by rationality and forward-induction reasoning (see Theorem 1). Recall that, for each  $i \in I$ ,  $T_i^k = \Theta_i \times M_i^k$  and  $T_{-i}^k = \Theta_{-i} \times M_{-i}^k$  are the spaces of  $(k)$  order personal features of  $i$  and  $-i$  respectively.

(Step 0) For all  $i \in I$ , let  $\mathbf{P}_i(0) = \Omega_i^k$ ,  $\mathbf{P}_{-i}(0) = \Omega_{-i}^k$  and  $\mathbf{P}(0) = \Omega^k$ .

(Step  $n > 0$ ) Assume that  $\mathbf{P}_i(m)$ ,  $\mathbf{P}_{-i}(m)$  and  $\mathbf{P}(m)$  have been defined for all  $m \in \{0, \dots, n-1\}$ . For each  $i \in I$ , let  $(z, \theta_i, \mu_i^k) \in \mathbf{P}_i(n)$  if there exists  $\mu_{i,k+1} \in M_{i,k+1}$  such that:

- (Coherence)  $(\mu_i^k, \mu_{i,k+1}) \in C_i^{k+1}$ ;
- (RP) for each  $h \in H$ ,  $\hat{\sigma}_i(\mu_{i,1})(r_{i,h}(\theta_i, \mu_{i,k+1}) | h) = 1$ ;
- (MC) for each  $h \in H$ , if  $h < z$ , then  $\hat{\sigma}_i(\mu_{i,1})(a_{i,h}(z) | h) > 0$ ;

<sup>36</sup> Indeed, own-plan independence is the key property (see Battigalli et al., 2019a). We consider the stronger OBI condition for expositional simplicity.

<sup>37</sup> See Section 4.2 in Battigalli and Dufwenberg (2019) and references therein.

<sup>38</sup> See Proposition 6 in the Appendix.

<sup>39</sup> Cf. Battigalli and Dufwenberg (2009), Sections 3–5.

- (Strong belief) for each  $m \in \{1, \dots, n - 1\}$ ,  $\mu_{i,k+1}$  strongly believes  $\mathbf{P}_{-i}(m)$ .

Finally, for all  $i \in I$ , let

$$\mathbf{P}_{-i}(n) = \bigcap_{j \in I \setminus \{i\}} \left( \mathbf{P}_j(n) \times \prod_{t \in I \setminus \{i,j\}} T_t^k \right) \text{ and } \mathbf{P}(n) = \bigcap_{i \in I} \left( \mathbf{P}_i(n) \times T_{-i}^k \right).$$

**Proposition 3.** For every  $n \in \mathbb{N}$ , the following are true:

- i) For every  $i \in I$ ,  $\mathbf{P}_i(n)$  is measurable.
- ii) For every  $i \in I$ ,  $\mathbf{P}_i(n) \subseteq \mathbf{P}_i(n - 1)$ ,  $\mathbf{P}_{-i}(n) \subseteq \mathbf{P}_{-i}(n - 1)$  and  $\mathbf{P}(n) \subseteq \mathbf{P}(n - 1)$ .

Therefore, the sequence of prediction sets  $(\mathbf{P}(n))_{n \in \mathbb{N}_0}$  is decreasing and it is standard to define  $\mathbf{P}(\infty) = \bigcap_{n \in \mathbb{N}_0} \mathbf{P}(n)$ . If

$\omega^k \in \mathbf{P}(\infty)$  then we say that  $\omega^k$  is **strongly rationalizable**.

If we replace condition MC above with *strict* material consistency, we obtain the decreasing sequence of events  $(\mathbf{P}^*(n))_{n \in \mathbb{N}_0}$  and its limit set  $\mathbf{P}^*(\infty)$ . Note that, in general, we cannot prove the non-emptiness of  $\mathbf{P}(\infty)$  or  $\mathbf{P}^*(\infty)$ . However, we can provide general and relevant sufficient conditions to obtain these results. Indeed, if we assume that  $(\Gamma, \nu) \in PG_k^{OBI}$ ,<sup>40</sup> then the behavioral equivalence and non-emptiness of  $\mathbf{P}(\infty)$  or  $\mathbf{P}^*(\infty)$  are obtained.<sup>41</sup>

**Proposition 4.** If  $(\Gamma, \nu) \in PG_k^{OBI}$ , then the following are true:

- i) For every  $n \in \mathbb{N} \cup \{\infty\}$ ,  $\text{proj}_Z \mathbf{P}(n) = \text{proj}_Z \mathbf{P}^*(n)$ .
- ii) For every  $n \in \mathbb{N} \cup \{\infty\}$ ,  $\mathbf{P}^*(n)$  is non-empty and compact.

The following example illustrates strong rationalizability in a psychological game where own-belief independence holds.

**Example 4.** Consider the game form with monetary payoffs of Fig. 2 with  $x = 1$ :  $S$  is a selfish action that maximizes the monetary payoff of Chloe. Suppose that Ann and Bob are (commonly known to be) selfish and risk neutral,<sup>42</sup> whereas Chloe may be affected by guilt aversion. Specifically, her belief-dependent utility is

$$u_c(z, \theta_c, \mu_{a,1}, \mu_{b,1}) = \pi_c(z) - \theta_c \sum_{j \in \{a,b\}} [\mathbb{E}_{\mu_{j,1}}(\pi_j | \emptyset) - \pi_j(z)]^+,$$

where  $\theta_c \in [\underline{\theta}_c, \bar{\theta}_c] \subseteq \mathbb{R}_+$  denotes her guilt sensitivity, that is, how much she dislikes to disappoint Ann and Bob. Since they obtain their maximal payoff after  $N$  and 0 after  $S$ , only the latter can disappoint them, with disappointment equal to the payoff they initially expected to get. Thus, the decision utilities of actions  $N$  and  $S$  given  $(D, L)$  are, respectively

$$\bar{u}_{c,(D,L)}(N, \theta_c, \mu_{c,2}) = 0$$

and

$$\begin{aligned} \bar{u}_{c,(D,L)}(S, \theta_c, \mu_{c,2}) &= 1 - \theta_c \mathbb{E}_{\mu_{c,2}} \left( \sum_{j \in \{a,c\}} \mathbb{E}_{\mu_{j,1}}(\pi_j | \emptyset) \middle| ((D, L), S) \right) \\ &= 1 - \theta_c \mathbb{E}_{\mu_{c,2}} \left( \sum_{j \in \{a,c\}} \mathbb{E}_{\mu_{j,1}}(\pi_j | \emptyset) \middle| (D, L) \right), \end{aligned}$$

where the latter equality holds by OAI. Note that Ann (Bob) can obtain \$2 for sure choosing  $U$  ( $R$ ). Therefore Ann (Bob) rationally chooses  $D$  ( $L$ ) only if s(he) expects to thereby get at least \$2, and chooses  $U$  ( $R$ ) otherwise. Since  $\max_{j \in \{a,b\}} \mathbb{E}_{\mu_{j,1}}(\pi_j | \emptyset) \leq 3$ , Chloe may rationally plan  $N$  only if  $1 - 6\theta_c \leq 0$ . If Chloe strongly believes in Ann’s and Bob’s

<sup>40</sup> Actually, it is sufficient to assume that  $i$ ’s utility does not depend on  $i$ ’s plan, which is just a (first-order) feature of  $i$ ’s belief hierarchy. We consider the stronger assumption in the text to simplify the exposition.

<sup>41</sup> The proof of Proposition 4 is available upon request. The first part follows from the fact that a behavior strategy is a sequential best reply to a conditional probability system (about the coplayers) if and only if every pure strategy in its “support” is also a sequential best reply. The proof of the second part adapts the proof of Theorem 13 in Battigalli and Dufwenberg (2009).

<sup>42</sup> That is,  $v_i(\pi(z), \theta, \mu) = \pi_i(z)$  for  $i = a, b$  and all  $z, \theta$ , and belief profiles  $\mu$ .

rationality, *ex ante* she cannot exclude any action pair of the coplayers, but *ex post* she would interpret  $(D, L)$  as signaling that  $\min_{j \in \{a,b\}} \mathbb{E}_{\mu_{j,1}}(\pi_j | \emptyset) \geq 2$ . Thus, we obtain two relevant thresholds for  $\theta_c$ ,  $1/4$  and  $1/6$ . We consider three relevant cases<sup>43</sup>:

► If  $\bar{\theta}_c < 1/6$  it is as if Chloe were unaffected by guilt aversion: rationality implies Chloe would choose selfishly upon observing  $(D, L)$  (1st step). Thus, Ann’s rationality and belief in others’ rationality implies that she chooses  $U$  (2nd step). Bob’s belief in rationality and in the coplayers’ belief in rationality implies that he expects to get \$3 by choosing  $L$ , which he rationally does (3rd step). With this, strong rationalizability implies  $(U, L)$ , which is correctly expected (4th step), on top of the expectation that Chloe would choose  $S$  (from the 2nd step). Thus, strong rationalizability yields  $(U, L)$ , the elimination process terminates in 4 steps, with the behavioral implications obtained in 3 steps.

► If  $\bar{\theta}_c > 1/4$ , Chloe may rationally plan to choose  $N$ ; if she also strongly believes in Ann’s and Bob’s rationality, she definitely plans to choose  $N$  (2nd step). Thus, Ann and Bob predict that Chloe would act non-selfishly if given the opportunity and Bob chooses  $L$  respectively (3rd step),<sup>44</sup> which is correctly expected, inducing Ann to rationally choose  $D$  (4th step), which is correctly expected (5th step). Thus, strong rationalizability yields  $((D, L), N)$ , the elimination process terminates in 5 steps, with the behavioral implications obtained in 4 steps.

► If  $1/6 < \bar{\theta}_c < 1/4$ , strong rationalizability has no behavioral implications: Chloe may rationally plan to choose  $N$ , but her strong belief in Ann’s and Bob’s rationality does not imply that she definitely plans to choose  $N$ . Thus, Ann and Bob may not trust her to act non-selfishly even if they believe that she is rational and that she strongly believes in their rationality. The elimination process changes slightly according to whether the commonly known upper bound  $\bar{\theta}_c$  is above or below  $1/4$ . If  $\bar{\theta}_c < 1/4$  the analysis is simpler, as rationality and mutual strong belief in rationality of any order do not have any implication for behavior, nor any joint implication for behavior and first-order beliefs beyond those of mere rationality. If instead  $\bar{\theta}_c > 1/4$ , Ann and Bob expect that high types of Chloe would choose  $N$  if given the opportunity (3rd step). Yet, they may also believe that high types have low probability; hence, such belief restriction has no impact on possible behaviors. ▲

### 7. Epistemic justification of strong rationalizability

In this section we provide an epistemic justification for the algorithm of strong rationalizability defined above. Recall that  $g_i : C_i^\infty \rightarrow \Delta_i^{H_i}(\Omega_{-i}^\infty)$  is the canonical homeomorphism of Proposition 2, and that  $\Omega_{-i}^\infty(h) = Z(h) \times T_{-i}^\infty$  for all  $h \in H$ .

**Definition 10.** The **strong belief operator** of  $i$  is a map  $SB_i : \mathcal{B}(\Omega_{-i}^\infty) \rightarrow 2^{\Omega_i^\infty}$  defined as

$$SB_i(E_{-i}) = \left\{ (z, \theta_i, \mu_i^\infty) \in \Omega_i^{\infty,*} : \forall h \in H, \Omega_{-i}^\infty(h) \cap E_{-i} \neq \emptyset \Rightarrow g_i(\mu_i^\infty)(E_{-i}|h) = 1 \right\},$$

for all  $E_{-i} \in \mathcal{B}(\Omega_{-i}^\infty)$ .

The following remark clarifies some of the properties of  $SB_i$ .

**Remark 6.** For all  $E_{-i} \in \mathcal{B}(\Omega_{-i}^\infty)$ ,  $SB_i(E_{-i})$  is measurable; if  $E_{-i}$  is closed, then  $SB_i(E_{-i})$  is closed.

We express our epistemic assumptions as events about behavior and personal features (personal traits and beliefs) of each player  $i$ , that is, measurable subsets of  $\Omega_i^\infty = Z \times T_i^\infty$ . For instance, we say that player  $i$  is rational and strongly believes the rationality of his coplayers at  $(z, t_i^\infty) \in \Omega_i^\infty$  if  $(z, t_i^\infty) \in R_i \cap SB_i(R_{-i}) \subseteq \Omega_i^\infty$ . The elements of  $R_i$  satisfy cross restrictions concerning actual behavior and the personal features  $t_i^\infty$  of player  $i$ , whereas the elements of  $SB_i(R_{-i})$  are just characterized by restrictions concerning the personal features of  $i$ , that is,  $\text{proj}_Z SB_i(R_{-i}) = Z$ . Similarly, the event that the coplayers  $-i$  are rational and strongly believe in their coplayers’ rationality is

$$\bigcap_{j \in I \setminus \{i\}} \left( (R_j \cap SB_j(R_{-j})) \times \prod_{\iota \in I \setminus \{i,j\}} T_\iota^\infty \right) \subseteq \Omega_{-i}^\infty.$$

Consider the following epistemic assumptions of increasing strength:

(Step 1) For every  $i \in I$ , let  $R_i(1) = R_i \subseteq \Omega_i^\infty$  and  $R_{-i}(1) = R_{-i} \subseteq \Omega_{-i}^\infty$ .

(Step  $n$ ) Assume that  $R_i(m)$  and  $R_{-i}(m)$  have been defined for every  $i \in I$  and  $m \in \{1, \dots, n-1\}$ , then define

$$R_i(n) = R_i(n-1) \cap SB_i(R_{-i}(n-1))$$

and

<sup>43</sup> In the appendix, we offer a complete formal analysis of each step  $\mathbf{P}(n)$ .

<sup>44</sup> Note that Bob expects to get \$3 with  $L$  independently of Ann’s choice.



$$R_{-i}(n) = \bigcap_{j \in I \setminus \{i\}} \left( R_j(n) \times \prod_{\ell \in I \setminus \{i, j\}} T_\ell^\infty \right).$$

With this, for every  $n \in \mathbb{N}$ , event  $R(n) = \bigcap_{i \in I} [R_i(n) \times T_i^\infty]$  represents the hypothesis of rationality (of all players) and  $n$ -mutual strong belief in rationality. By definition, it follows that, for every  $i \in I$  and  $n \in \mathbb{N}$ ,  $R_i(n+1) \subseteq R_i(n)$  and  $R(n+1) \subseteq R(n)$ . Therefore, the event in  $\Omega^\infty$  that represents *rationality and common strong belief in rationality* is  $R(\infty) = \bigcap_{n \in \mathbb{N}} R(n)$ . As explained in Battigalli and Siniscalchi (2002) these assumptions require, on top of rationality, that each player ascribes to his coplayers the highest degree of “strategic sophistication” consistent with their observed behavior, that is, given any  $h \in H$  player  $i$  assigns probability 1 to  $R_{-i}(n_{-i,h}^*)$ , where

$$n_{-i,h}^* = \max \{n \in \mathbb{N} : \Omega_{-i}^\infty(h) \cap R_{-i}(n) \neq \emptyset\}.$$

We can now state the main result of this paper.

**Theorem 1.** For every  $n \in \mathbb{N}$ ,

$$\forall i \in I, \mathbf{P}_i(n) = \text{proj}_{\Omega_i^k} R_i(n) \quad \text{and} \quad \mathbf{P}(n) = \text{proj}_{\Omega^k} R(n).$$

The theorem characterizes the utility-relevant implications of rationality and  $n$ -mutual strong belief in rationality, thus providing an epistemic justification for the strong rationalizability algorithm in  $k$ th-order psychological games. The proof is in the Appendix, here we give a brief sketch. The harder part of the proof is to show that, for all  $n \in \mathbb{N}$  and for every  $i \in I$ , given an element  $(z, \theta_i, \mu_i^k) \in \mathbf{P}_i(n)$  we can find an infinite hierarchy  $\mu_i^\infty$  of conditional beliefs for  $i$  which is consistent with  $\mu_i^k$  and yields a personal state in  $R_i(n)$ . Here, the factorization in Eq. (3) allows us to invoke Lemma 3 to show the existence of a consistent ICPS defined over the space of total uncertainty for  $i$ . Then, we use the canonical homeomorphism  $g_i$  (see Proposition 2) to derive  $\mu_i^\infty$ .

Relying on Theorem 1, under the assumption that the psychological utility  $v_i$  of every player  $i$  does not depend on his own beliefs, one can use standard compact-continuity arguments to prove the following result.

**Proposition 5.** If  $(\Gamma, v) \in PG_k^{OBI}$ , then

$$\mathbf{P}^*(\infty) = \text{proj}_{\Omega^k} R^*(\infty),$$

where  $\mathbf{P}^*(\infty)$  is non-empty and compact.

Under OBI we can rely on the properties of state-dependent expected utility. As for standard games, we can focus without loss of generality on deterministic rational plans (see Proposition 4), obtain a decreasing sequence of compact events and apply the finite intersection property to show that the limit set obtained from our algorithm is non-empty and characterizes the utility-relevant implications of rationality and common strong belief in rationality.

## 8. Discussion

We discuss some features of our approach and simplifying assumptions, hinting at extensions and generalizations.

**External uncertainty and behavior** Players are uncertain about the mental states (beliefs) of coplayers, their traits, and behavior. First-order beliefs are beliefs about the non-mental, or external aspects of what is uncertain. In particular, beliefs about behavior are (aspects of) first-order beliefs. The representation of behavior is uncontroversial for simultaneous-move games, where it is given by the profile of actions simultaneously chosen by the players. In dynamic games, instead, there are multiple ways to model the space of possible behaviors. The more traditional one is to represent behavior as a conjunction of subjunctive conditional statements of the form “if  $h$  were reached,  $j$  would take action  $a_j$ .” Formally, the conjunction of all such statements for player  $j$  is a pure strategy. Thus, the traditional approach is to represent behavior as a profile of pure strategies. Yet, if we interpret strategies as plans in the minds of the players, they should belong to the mental part of the state of the world, not the external one. Modeling strategies as plans in the minds of players is essential for many applications of psychological game theory.<sup>45</sup> Thus, if we follow the traditional representation of behavior, we have to allow for two mathematical objects that can legitimately be called “strategy”: the objective description of behavioral subjunctive conditionals, and the subjective plans in the minds of players.<sup>46</sup> This may create confusion and misunderstandings. Therefore, we follow Battigalli et al. (2013)<sup>47</sup> and represent behavior as the actual sequence of actions profiles chosen by the players,

<sup>45</sup> See Battigalli and Dufwenberg (2019) and Battigalli et al. (2019a).

<sup>46</sup> See, for example, Section 5 of Battigalli and Siniscalchi (1999).

<sup>47</sup> As well as Battigalli and Dufwenberg (2019), and Battigalli et al. (2019a, 2019b).

whereas strategies are expressed as beliefs about own behavior and belong to the mental part of the state of the world. This approach, however, has costs as well. Strategic reasoning, in essence, requires players to have beliefs about how coplayers would react to what they observe, i.e., beliefs about behavioral subjunctive conditionals. Yet, such conditional are not part of our formal language. We circumvent this difficulty by *expressing the probabilities of conditionals as conditional probabilities*, as in Kuhn’s (1953) transformation from mixed to behavior strategies. For example, if  $i$  at  $h$  believes with probability  $p$  that  $j$  would take action  $a_j$  should he choose action  $a_i$ , then we ascribe to  $i$  a first-order belief  $\mu_i^1$  such that  $\mu_i^1(a_j|h, a_i) = p$ . In sum, we do not claim that the present approach dominates the traditional one, but we find it germane to the representation of psychological games where players’ intentions, hence their subjective plans, are key.

**Imperfectly observable actions and cognitive rationality** We assumed that the actions of previous stages are perfectly observed because this allows us to simplify the notation and streamline the analysis. But we can prove our results for all multi-stage games. Here we sketch this generalization. In doing so, we highlight the possibility and convenience to consider more general belief spaces by dropping the necessity of other cognitive rationality properties besides the coherence of belief hierarchies. Let  $Y_i$  denote a set of “personal outcomes” or “messages” that player  $i$  may observe as the play unfolds. For each stage  $t$  such that  $A^t \cap \bar{H} \neq \emptyset$  there is a feedback function  $f_{i,t} : A^t \cap \bar{H} \rightarrow Y_i$  that represents the *flow* of information acquired by  $i$  at the end of stage  $t$  according to the rules of the game, such as his stage-game monetary payoff, or his cumulated monetary payoff in a repeated game. If  $i$  has perfect recall, at the end of stage  $t$  (and the beginning of stage  $t + 1$ ) he remembers the personal history  $h_i^t = (a_{i,k}, y_{i,k})_{k=1}^t$  of actions and personal outcomes and can back out the **information set**  $[h_i^t] \subseteq \bar{H}$  of histories consistent with it.<sup>48</sup> Let  $\mathbf{H}_i$  denote the collection containing the information sets  $[h_i^t]$  as well as similarly defined “interim information sets”  $[h_i^t, a_{i,t+1}]$ . If the player is cognitively rational, he remembers all previous choices and signals; thus, as the play unfolds, he conditions on the *stock* of cumulated information represented by elements of  $\mathbf{H}_i$ . If he is not fully cognitively rational and forgets his previous choices and signals, he only conditions on the last piece of information  $y_{i,t}$  he just obtained, that is, the set  $f_{i,t}^{-1}(y_{i,t}) \subseteq \bar{H}$ . Thus, it makes sense to define the spaces of conditional beliefs so that conditioning refers only to “memoryless” information sets of the form  $f_{i,t}^{-1}(y_{i,t})$ , and possibly without even assuming the chain rule, whose bite depends heavily on (perfect) memory, which makes the collection of conditioning events ordered by set inclusion a tree.<sup>49</sup> The key observation is that as long as we define (cognitive) *rationality* so as to encompass perfect recall, application of the chain rule, and coherence belief hierarchies, our analysis can be extended to this more general environment.

**Generalized psychological utilities** In our analysis, the decision utility of action  $a_i$  given history  $h$  is the subjective expectation of psychological utility  $v_i$  conditional on  $(h, a_i)$  given personal traits  $\theta_i$  and  $k + 1$  order belief  $\mu_{i,k+1}$ . With this, we obtain a continuous “local” utility function  $\bar{u}_{i,h}(a_i, \theta_i, \mu_{i,k+1})$ . We emphasized that only the local utility functions  $(\bar{u}_{i,h})_{(i,h) \in I \times H}$  matter for our epistemic analysis. Thus, the results in this article are valid for psychological games with more general forms of psychological preferences. For example, we may obtain  $\bar{u}_{i,h}$  as a local “distortion” of the conditional expectation of “experience utility”  $v_i$  (see Battigalli et al., 2019a and 2019b).

**Infinite games** We considered finite games forms, but our analysis extends to a large class of multi-stage games where players’ feasible action sets are finite at all histories of height 2 or more, and their psychological utility functions do not depend on terminal beliefs, as we can still adapt the techniques of Battigalli and Tebaldi (2019) to such games. This covers, for example, all compact-continuous games with simultaneous moves where utility depends only on initial beliefs.

**Rationalizable self-confirming equilibrium** The self-confirming equilibrium (SCE) concept — a generalization of Nash equilibrium — characterizes the steady states of learning dynamics in games played recurrently. According to SCE, agents are asymptotic empiricists who best reply to confirmed, but possibly false beliefs. In particular, players need not believe in the strategic sophistication of others. In a (strongly) rationalizable SCE (RSCE) agents (strongly) believe in the strategic sophistication of others. We can provide an algorithmic definition of RSCE for multi-stage ( $k$  order psychological) games and an epistemic justification of RSCE.

**Restrictions on low-order beliefs and  $\Delta$ -rationalizability** In applications, it is often natural to impose some restrictions on low-order beliefs and assume that such restrictions shape strategic reasoning (see, e.g., Battigalli and Tebaldi, 2019 and references therein). This yields a modified notion of strong rationalizability, called strong  $\Delta$ -rationalizability (where  $\Delta$

<sup>48</sup> That is,  $[h_i^t] = \left\{ (a'_{i,k})_{k=1}^t \in \bar{H} : \forall k, a'_{i,k} = a_{i,k}, f_{i,k}(a'_{i,1}, \dots, a'_{i,k}) = y_{i,k} \right\}$ . By construction, this collection of information sets satisfies the standard perfect-recall conditions.

<sup>49</sup> As we consider larger belief spaces allowing for more inconsistencies, it becomes more important to specify how psychological utility may depend on beliefs. The structural assumption should be that utility depends only on realized beliefs about outcomes, or behavior, or traits, and higher-order realized beliefs. This is important also for rational players, whose beliefs are fully consistent, because their utility may depend on the unknown beliefs of possibly irrational coplayers. See Battigalli et al. (2019a).

represents the restricted set of profiles of beliefs). We can define strong  $\Delta$ -rationalizability for  $k$  order psychological games and provide an epistemic justification of this solution concept.

**Iterated conditional dominance** We defined strong rationalizability for  $k$ -th order psychological games and provided an epistemic foundation for this solution procedure, thus extending the results proved by Battigalli and Siniscalchi (2002) for finite standard games, and by Battigalli and Tebaldi (2019) for a class of infinite standard games. Shimoji and Watson (1998) proved that strong rationalizability in finite standard games can be algorithmically implemented by a procedure of iterated elimination of conditionally dominated strategies. We can extend their characterization result to the class of infinite games analyzed by Battigalli and Tebaldi (2019) and to the  $k$ -th order psychological games that satisfy own-belief independence (see Battigalli and Corrao, 2019). The result does not hold for general psychological games (see Mourmans, 2019).

## 9. Appendix

### 9.1. Complete analysis of strong rationalizability in Example 4

Recall that  $\Theta$  is isomorphic to  $[\underline{\theta}_c, \bar{\theta}_c]$ , that is, it is common knowledge that  $\theta_a = \theta_b = 0$ , and  $\underline{\theta}_c \leq \theta_c \leq \bar{\theta}_c$ . The utility functions are  $u_i(z, \theta, \mu_1) = \pi_i(z)$  for  $i \in \{a, b\}$ , and

$$u_c(z, \theta_c, \mu_{a,1}, \mu_{b,1}) = \pi_c(z) - \theta_c \sum_{j \in \{a,b\}} [\mathbb{E}_{\mu_{j,1}}(\pi_j | \emptyset) - \pi_j(z)]^+.$$

Hence,

$$\begin{aligned} \mathbf{P}_a(1) &= \{(z, \mu_{a,1}) : \alpha_a(z) = U, \hat{\sigma}_a(\mu_{a,1})(U | \emptyset) > 0, \mathbb{E}_{\mu_{a,1}}(\pi_a | \emptyset) = 2, \mathbb{E}_{\mu_{a,1}}(\pi_a | D) \leq 2\} \\ &\cup \{(z, \mu_{a,1}) : \alpha_a(z) = D, \hat{\sigma}_a(\mu_{a,1})(D | \emptyset) > 0, \mathbb{E}_{\mu_{a,1}}(\pi_a | \emptyset) = \mathbb{E}_{\mu_{a,1}}(\pi_a | D) \geq 2\}, \\ \mathbf{P}_b(1) &= \{(z, \mu_{b,1}) : \alpha_b(z) = R, \hat{\sigma}_b(\mu_{b,1})(R | \emptyset) > 0, \mathbb{E}_{\mu_{b,1}}(\pi_b | \emptyset) = 2, \mathbb{E}_{\mu_{b,1}}(\pi_b | L) \leq 2\} \\ &\cup \{(z, \mu_{b,1}) : \alpha_b(z) = L, \hat{\sigma}_b(\mu_{b,1})(L | \emptyset) > 0, \mathbb{E}_{\mu_{b,1}}(\pi_b | \emptyset) = \mathbb{E}_{\mu_{b,1}}(\pi_b | L) \geq 2\}, \end{aligned}$$

because  $a$  and  $b$  can secure payoff 2 by choosing  $U$  and  $R$  respectively.

As for  $c$ , there are two thresholds:  $\theta_c = 1/6$  makes  $c$  indifferent if she believes  $a$  and  $b$  initially expected the maximum payoff 3, while  $\theta_c = 1/4$  makes  $c$  indifferent if she thinks that  $a$  and  $c$  expected the payoff they can secure, 2. With this, we consider 3 cases:

- $\bar{\theta}_c < 1/6$  is strategically equivalent to common knowledge that also  $c$  is selfish. Thus, mere rationality rules out the non-selfish action  $N$ , and strong rationalizability yields  $(U, L)$  with the expectation that Chloe would choose  $S$ :

$$\begin{aligned} \mathbf{P}_c(1) &= \mathbf{P}_c(2) = \{(z, \theta_c, \mu_{1,c}) : (\alpha_a(z), \alpha_b(z)) \neq (D, L), \hat{\sigma}_c(S | (D, L)) = 1\} \\ &\cup \{(z, \theta_c, \mu_{1,c}) : z = ((D, L), S), \hat{\sigma}_c(S | (D, L)) = 1\}, \\ \mathbf{P}_a(2) &= \{(z, \mu_{a,1}) : \alpha_a(z) = U, \hat{\sigma}_a(\mu_{a,1})(U | \emptyset) = 1, \mathbb{E}_{\mu_{a,1}}(\pi_a | D) = 0\}, \\ \mathbf{P}_b(2) &= \{(z, \mu_{b,1}) \in \mathbf{P}_b(1) : \mathbb{E}_{\mu_{b,1}}(\pi_b | (D, L)) = 0\} \\ \mathbf{P}_a(3) &= \mathbf{P}_a(2), \\ \mathbf{P}_b(3) &= \{(z, \mu_{b,1}) \in \mathbf{P}_b(2) : \alpha_b(z) = L, \mathbb{E}_{\mu_{b,1}}(\pi_b | \emptyset) = \mathbb{E}_{\mu_{a,1}}(\pi_a | L) = 3\}, \\ \mathbf{P}_c(3) &= \{(z, \theta_c, \mu_{1,c}) \in \mathbf{P}_c(2) : \mu_{c,1}(U | \emptyset) = 1\}, \\ \mathbf{P}_a(4) &= \{(z, \mu_{a,1}) \in \mathbf{P}(3) : \mu_{a,1}(L | \emptyset) = 1\}, \\ \mathbf{P}_b(4) &= \mathbf{P}_b(3), \\ \mathbf{P}_c(4) &= \{(z, \theta_c, \mu_{1,c}) \in \mathbf{P}_c(3) : \mu_{c,1}(L | \emptyset) = 1\}, \\ \forall n > 4, \forall i \in I, \mathbf{P}_i(n) &= \mathbf{P}_i(4). \end{aligned}$$

- $\underline{\theta}_c > 1/4$  means it is common knowledge that  $c$  is so averse to guilt that if she strongly believes in  $a$ 's and  $c$ 's rationality, then she wants to take the non-selfish action upon observing  $(D, L)$ . Thus, while mere rationality does not rule out any plan of  $c$  (because  $\underline{\theta}_c > 1/6$ ) and just imposes material consistency, strong rationalizability yields  $(D, L)$  with the correct expectation that  $c$  chooses  $N$ :

$$\begin{aligned} \mathbf{P}_c(1) &= \mathbf{P}_c^{\text{high.G}}(1) := \{(z, \theta_c, \mu_{1,c}) : (\alpha_a(z), \alpha_b(z)) \neq (D, L)\} \\ &\cup \{(z, \theta_c, \mu_{1,c}) : (\alpha_a(z), \alpha_b(z)) = (D, L), \hat{\sigma}_c(\alpha_c(z) | (D, L)) > 0\}, \end{aligned}$$

$$\begin{aligned}
 \mathbf{P}_a(2) &= \mathbf{P}_a(1), \\
 \mathbf{P}_b(2) &= \mathbf{P}_b(1), \\
 \mathbf{P}_c(2) &= \{(z, \theta_c, \mu_{1,c}) \in \mathbf{P}_c(1) : \hat{\sigma}_c(N|(D, L)) = 1\}, \\
 \mathbf{P}_a(3) &= \{(z, \mu_{a,1}) \in \mathbf{P}_a(2) : \mathbb{E}_{\mu_{a,1}}(\pi_a|(D, L)) = 3\}, \\
 \mathbf{P}_b(3) &= \{(z, \mu_{b,1}) \in \mathbf{P}_b(2) : \alpha_b(z) = L, \mathbb{E}_{\mu_{b,1}}(\pi_b|L) = 3\}, \\
 \mathbf{P}_c(3) &= \mathbf{P}_c(2), \\
 \mathbf{P}_a(4) &= \{(z, \mu_{a,1}) \in \mathbf{P}_a(3) : \alpha_a(z) = D, \mathbb{E}_{\mu_{a,1}}(\pi_a|D) = 3\}, \\
 \mathbf{P}_b(4) &= \mathbf{P}_b(3), \\
 \mathbf{P}_c(4) &= \{(z, \theta_c, \mu_{1,c}) \in \mathbf{P}_c(3) : \mu_{1,c}(L|\emptyset) = 1\}, \\
 \mathbf{P}_a(5) &= \mathbf{P}_a(4), \\
 \mathbf{P}_b(5) &= \{(z, \mu_{b,1}) \in \mathbf{P}_b(4) : \mu_{b,1}(D|\emptyset) = 1\}, \\
 \mathbf{P}_c(5) &= \{(z, \theta_c, \mu_{1,c}) \in \mathbf{P}_c(4) : \mu_{1,c}(D|\emptyset) = 1\}, \\
 \forall n > 4, \forall i \in I, \mathbf{P}_i(n) &= \mathbf{P}_i(5).
 \end{aligned}$$

- $1/6 < \underline{\theta}_c < 1/4$  means that  $a$  and  $b$  do not know enough to predict  $c$ 's choice; thus, neither rationality, nor strong rationalizability rule out any behavior. There are two relevant sub-cases:
  - If  $1/6 < \underline{\theta}_c \leq \bar{\theta}_c < 1/4$ , then strong belief in rationality is not enough to compel any type of  $c$  to take the non selfish action  $N$  upon observing  $(D, L)$ . Thus, the first step is as in the previous case, but unlike the previous case there are no further restrictions in the following steps:

$$\begin{aligned}
 \mathbf{P}_c(1) &= \mathbf{P}_c^{\text{high.G}}(1), \\
 \forall n > 1, \forall i \in I, \mathbf{P}_i(n) &= \mathbf{P}_i(1).
 \end{aligned}$$

- If  $1/6 < \underline{\theta}_c < 1/4 < \bar{\theta}_c$ , strong rationalizability rules out some first-order beliefs about  $c$  because high-guilt types would act non-selfishly upon observing  $(D, L)$ :

$$\begin{aligned}
 \mathbf{P}_c(1) &= \mathbf{P}_c^{\text{high.G}}(1), \\
 \mathbf{P}_a(2) &= \mathbf{P}_a(1), \\
 \mathbf{P}_b(2) &= \mathbf{P}_b(1), \\
 \mathbf{P}_c(2) &= \left\{ (z, \theta_c, \mu_{1,c}) \in \mathbf{P}_c(1) : \theta_c > \frac{1}{4}, \hat{\sigma}_c(N|(D, L)) = 1 \right\} \\
 &\quad \cup \left\{ (z, \theta_c, \mu_{1,c}) \in \mathbf{P}_c(1) : \theta_c \leq \frac{1}{4} \right\}, \\
 \mathbf{P}_i(3) &= \left\{ (z, \mu_{i,1}) \in \mathbf{P}_i(2) : \mu_{i,1} \left( \left\{ (z, \theta_c) : z = ((D, L), S), \theta_c > \frac{1}{4} \right\} \middle| (D, L) \right) = 0 \right\} \quad (i \in \{a, b\}), \\
 \mathbf{P}_c(3) &= \mathbf{P}_c(2) \\
 \forall n > 3, \forall i \in I, \mathbf{P}_i(n) &= \mathbf{P}_i(3).
 \end{aligned}$$

### 9.2. Proofs of the main results

We start by proving a general result that implies Proposition 1. For all measurable functions  $\zeta_{-i} : T_{-i} \rightarrow \Sigma_{-i}$  and probability measures  $\eta_i \in \Delta(T_{-i})$  we define  $\zeta_{-i}^{\eta_i} \in \Sigma_{-i}$  as

$$\zeta_{-i}^{\eta_i}(a_{-i}|h) = \int_{T_{-i}} \zeta_{-i}^{t_{-i}}(a_{-i}|h) \eta_i(dt_{-i})$$

for all  $h \in H$  and  $a_{-i} \in A_{-i}(h)$ . Also, recall that for all  $(\sigma_i, \sigma_{-i}) \in \Sigma_i \times \Sigma_{-i}$ , and  $h_i, h'_i \in H_i$  with  $h'_i \geq h_i$ , we let

$$\mathbb{P}_{\sigma_i, \sigma_{-i}}(h'_i|h_i)$$

denote the probability of  $h'_i$  given  $h_i$  induced by  $(\sigma_i, \sigma_{-i})$ .

**Proposition 6.** Fix any  $\mu_i \in [\Delta(Z \times T_{-i})]^{H_i}$ . The following are equivalent:

- i)  $\mu_i$  is a CPS that satisfies OAI;
- ii)  $\mu_i$  is a CPS such that for all  $h \in H$ ,  $(a_i, a_{-i}) \in A(h)$ , and measurable  $E_{-i} \subseteq T_{-i}$ ,

$$\mu_i(Z(h, a_{-i}) \times E_{-i} | h) = \mu_i(Z(h, (a_i, a_{-i})) \times E_{-i} | h, a_i);$$

- iii) there exist  $\sigma_i \in \Sigma_i$ , a measurable function  $\zeta_{-i} : T_{-i} \rightarrow \Sigma_{-i}$ , and a vector of probability measures  $\eta_i \in [\Delta(T_{-i})]^H$  such that, for all  $h, h' \in H$  with  $h' \geq h$ ,  $z \in Z$ , and measurable  $E_{-i} \subseteq T_{-i}$ ,

$$\mathbb{P}_{\sigma_i, \zeta_{-i}}^{\eta_i}(h' | h) > 0 \implies \eta_{i, h'}(E_{-i}) = \frac{1}{\mathbb{P}_{\sigma_i, \zeta_{-i}}^{\eta_i}(h' | h)} \int_{E_{-i}} \mathbb{P}_{\sigma_i, \zeta_{-i}}^{\eta_i}(h' | h) \eta_{i, h}(dt_{-i}) \tag{6}$$

and

$$\begin{aligned} \mu_i(\{z\} \times E_{-i} | h) &= \int_{E_{-i}} (\sigma_i, \zeta_{-i}^{t_{-i}})(z | h) \eta_{i, h}(dt_{-i}), \\ \mu_i(\{z\} \times E_{-i} | h, a_i) &= \int_{E_{-i}} (\sigma_i, \zeta_{-i}^{t_{-i}})(z | h, a_i) \eta_{i, h}(dt_{-i}), \end{aligned}$$

The pair  $(\sigma_i, \eta_i)$  is unique and, given  $\eta_i$ , each function  $t_{-i} \mapsto \zeta_{-i}^{t_{-i}}(\cdot | h)$  is  $\eta_{i, h}$ -a.e. uniquely defined.

**Proof of Proposition 6.** i)  $\implies$  ii) Fix  $h \in H$ ,  $(a_i, a_{-i}) \in A(h)$ , and a measurable  $E_{-i} \subseteq T_{-i}$ . With this,

$$\begin{aligned} \mu_i(Z(h, a_{-i}) \times E_{-i} | h) &= \sum_{a'_i \in A_i(h)} \mu_i(Z(h, (a'_i, a_{-i})) \times E_{-i} | h) \\ &= \sum_{a'_i \in A_i(h)} \mu_i(Z(h, (a'_i, a_{-i})) \times E_{-i} | h, a'_i) \mu_i(a'_i | h) \\ &= \mu_i(Z(h, (a_i, a_{-i})) \times E_{-i} | h, a_i) \left( \sum_{a'_i \in A_i(h)} \mu_i(a'_i | h) \right) \\ &= \mu_i(Z(h, (a_i, a_{-i})) \times E_{-i} | h, a_i), \end{aligned}$$

where the third equality holds by OAI.

ii)  $\implies$  iii) For all  $h \in H$ ,  $a \in A(h)$ ,  $t_{-i} \in T_{-i}$ , and measurable  $E_{-i} \subseteq T_{-i}$ , define

$$\begin{aligned} \sigma_i(a_i | h) &= \mu_i(a_i | h), \\ \zeta_{-i}^{t_{-i}}(a_{-i} | h) &= \mu_i(Z(h, a_{-i}) | h, t_{-i}), \\ \eta_i(E_{-i} | h) &= \mu_i(Z \times E_{-i} | h), \end{aligned}$$

where  $t_{-i} \mapsto \mu_i(Z(h, a_{-i}) | h, t_{-i})$  is a version of the conditional probability of  $Z(h, a_{-i})$  given the probability measure  $\mu_i(\cdot | h) \in \Delta(Z \times T_{-i})$ . Note that, by construction, we have

$$\int_{T_{-i}} \zeta_{-i}^{t_{-i}}(a_{-i} | h) \eta_{i, h}(dt_{-i}) = \zeta_{-i}^{\eta_i}(a_{-i} | h) = \hat{\sigma}_{-i}(a_{-i} | h).$$

We now prove the following claim: For all  $h \in H$ ,  $a \in A(h)$ ,  $z \in Z$ , and measurable  $E_{-i} \subseteq T_{-i}$ ,

$$\mu_i(\{z\} \times E_{-i} | h) = \int_{E_{-i}} \mathbb{P}_{\sigma_i, \zeta_{-i}}^{\eta_i}(z | h) \eta_{i, h}(dt_{-i})$$

and

$$\mathbb{P}_{\sigma_i, \zeta_{-i}}^{\eta_i}((h, a) | h) > 0 \implies \eta_{i, (h, a)}(E_{-i}) = \frac{1}{\mathbb{P}_{\sigma_i, \zeta_{-i}}^{\eta_i}((h, a) | h)} \int_{E_{-i}} \mathbb{P}_{\sigma_i, \zeta_{-i}}^{\eta_i}((h, a) | h) \eta_{i, h}(dt_{-i}).$$

We prove this claim by induction on the height  $L(h) = \max\{\ell(z) - \ell(h) : z \in Z(h)\}$  of histories. Pick any  $h \in \bar{H}$  such that  $L(h) = 1$ , that is,  $h$  is preterminal. Then, for all  $a \in A(h)$  and measurable  $E_{-i} \subseteq T_{-i}$ ,

$$\begin{aligned} \mu_i(\{(h, a)\} \times E_{-i}|h) &= \mu_i(\{(h, a_i)\} \times T_{-i}|h) \mu_i(\{(h, a)\} \times E_{-i}|h, a_i) \\ &= \mu_i(\{(h, a_i)\} \times T_{-i}|h) \mu_i(\{(h, a_{-i})\} \times E_{-i}|h) \\ &= \int_{E_{-i}} \mu_i(\{(h, a_i)\} \times T_{-i}|h) \mu_i(\{(h, a_{-i})\} |h, t_{-i}) \text{marg}_{T_{-i}} \mu_{i,h} (dt_{-i}) \\ &= \int_{E_{-i}} \mathbb{P}_{\sigma_i, \zeta_{-i}^{t_{-i}}}(z|h) \eta_{i,h} (dt_{-i}). \end{aligned}$$

Next, pick any  $a \in A(h)$  such that  $\mathbb{P}_{\sigma_i, \zeta_{-i}^{\eta_i}}((h, a)|h) > 0$ , that is,  $\mu_i(\{(h, a)\} \times T_{-i}|h) > 0$ . Then

$$\begin{aligned} \eta_{i,(h,a)}(E_{-i}) &= \mu_i(Z \times E_{-i}|(h, a)) \\ &= \mu_i(\{(h, a)\} \times E_{-i}|(h, a)) \\ &= \frac{1}{\mu_i(\{(h, a)\} \times T_{-i}|h)} \mu_i(\{(h, a)\} \times E_{-i}|h) \\ &= \frac{1}{\mathbb{P}_{\sigma_i, \zeta_{-i}^{\eta_i}}((h, a)|h)} \int_{E_{-i}} \mathbb{P}_{\sigma_i, \zeta_{-i}^{t_{-i}}}(z|h) \eta_{i,h} (dt_{-i}). \end{aligned}$$

Suppose by way of induction that the claim holds for every history  $h' \in H$  with  $L(h') \leq n$ . Pick any  $h \in \bar{H}$  such that  $L(h) = n + 1$ . Then, for all  $z \in Z(h)$  and measurable  $E_{-i} \subseteq T_{-i}$ ,

$$\begin{aligned} \mu_i(\{z\} \times E_{-i}|h) &= \mu_i(Z((h, a_{h,z})) \times T_{-i}|h) \mu_i(\{z\} \times E_{-i}|h, a_{h,z}) \\ &= \mu_i(Z((h, a_{i,h,z})) \times T_{-i}|h) \int_{E_{-i}} \mu_i(Z((h, a_{-i,h,z})) |h, t_{-i}) \mathbb{P}_{\sigma_i, \zeta_{-i}^{t_{-i}}}(z|h, a_{h,z}) \eta_{i,(h,a_{h,z})} (dt_{-i}) \\ &= \int_{E_{-i}} \sigma_i(a_{i,h,z}|h) \zeta_{-i}^{t_{-i}}(a_{-i,h,z}|h) \mathbb{P}_{\sigma_i, \zeta_{-i}^{t_{-i}}}(z|h, a_{h,z}) \eta_{i,(h,a_{h,z})} (dt_{-i}) \\ &= \int_{E_{-i}} \mathbb{P}_{\sigma_i, \zeta_{-i}^{t_{-i}}}(z|h) \eta_{i,h} (dt_{-i}), \end{aligned}$$

where  $a_{h,z}$  is the unique profile of feasible actions at  $h$  implied by  $z$  and  $a_{i,h,z}$  and  $a_{-i,h,z}$  are similarly defined. Next, consider some  $a \in A(h)$  such that  $\mathbb{P}_{\sigma_i, \zeta_{-i}^{\eta_i}}((h, a)|h) > 0$ , that is,  $\mu_i(\{(h, a)\} \times T_{-i}|h) > 0$ . We have

$$\begin{aligned} \eta_{i,(h,a)}(E_{-i}) &= \mu_i(Z \times E_{-i}|(h, a)) \\ &= \mu_i(Z(h, a) \times E_{-i}|(h, a)) \\ &= \frac{1}{\mu_i(Z(h, a) \times T_{-i}|h)} \mu_i(Z(h, a) \times E_{-i}|h) \\ &= \frac{1}{\mathbb{P}_{\sigma_i, \zeta_{-i}^{\eta_i}}((h, a)|h)} \sum_{z \in Z(h,a)} \int_{E_{-i}} \mathbb{P}_{\sigma_i, \zeta_{-i}^{t_{-i}}}(z|h) \eta_{i,h} (dt_{-i}) \\ &= \frac{1}{\mathbb{P}_{\sigma_i, \zeta_{-i}^{\eta_i}}((h, a)|h)} \int_{E_{-i}} \sum_{z \in Z(h,a)} \left[ \mathbb{P}_{\sigma_i, \zeta_{-i}^{t_{-i}}}(z|h) \right] \eta_{i,h} (dt_{-i}) \\ &= \frac{1}{\mathbb{P}_{\sigma_i, \zeta_{-i}^{\eta_i}}((h, a)|h)} \int_{E_{-i}} \mathbb{P}_{\sigma_i, \zeta_{-i}^{t_{-i}}}(z|h) \eta_{i,h} (dt_{-i}), \end{aligned}$$

proving the statement. Note that an analogous argument shows that, for all  $h, h' \in H$  with  $h' \geq h$ ,  $a_i \in A_i(h)$ ,  $z \in Z$  and measurable  $E_{-i} \subseteq T_{-i}$ ,

$$\mu_i(\{z\} \times E_{-i}|h, a_i) = \int_{E_{-i}} (\sigma_i, \zeta_{-i}^{t_{-i}})(z|h, a_i) \eta_{i,h} (dt_{-i})$$

and

$$\mathbb{P}_{\sigma_i, \zeta_{-i}}^{\eta_{i,h}}(h'|h) > 0 \implies \eta_{i,h'}(E_{-i}) = \frac{1}{\mathbb{P}_{\sigma_i, \zeta_{-i}}^{\eta_{i,h}}(h'|h)} \int_{E_{-i}} \mathbb{P}_{\sigma_i, \zeta_{-i}}^{\eta_{i,h}}(h'|h) \eta_{i,h}(\mathbf{dt}_{-i}).$$

iii)  $\implies$  i) Consider a vector of probability measures  $\mu_i \in [\Delta(Z \times T_{-i})]^{H_i}$  satisfying the factorization  $(\sigma_i, \zeta_{-i}, \eta_i)$  defined in point iii). We need to show that  $\mu_i$  is a ICPS. Fix  $h_i \in H$  and note that

$$\begin{aligned} \mu_i(Z(h_i) \times T_{-i}|h_i) &= \sum_{z \in Z(h_i)} \int_{T_{-i}} \mathbb{P}_{\sigma_i, \zeta_{-i}}^{\eta_i}(z|h_i) \eta_{i,h}(\mathbf{dt}_{-i}) \\ &= \int_{T_{-i}} \left[ \sum_{z \in Z(h_i)} \mathbb{P}_{\sigma_i, \zeta_{-i}}^{\eta_i}(z|h_i) \right] \eta_{i,h}(\mathbf{dt}_{-i}) \\ &= \int_{T_{-i}} \eta_{i,h}(\mathbf{dt}_{-i}) = 1. \end{aligned}$$

Next, fix  $h, h' \in H$ ,  $z \in Z$  such that  $z \geq h' \geq h$ , and a measurable set  $E_{-i} \subseteq T_{-i}$ . We have

$$\begin{aligned} \mu_i(\{z\} \times E_{-i}|h) &= \int_{E_{-i}} \mathbb{P}_{\sigma_i, \zeta_{-i}}^{\eta_{i,h}}(z|h) \eta_{i,h}(\mathbf{dt}_{-i}) \\ &= \mathbb{P}_{\sigma_i, \zeta_{-i}}^{\eta_{i,h}}(z|h) \frac{\int_{E_{-i}} \mathbb{P}_{\sigma_i, \zeta_{-i}}^{\eta_{i,h'}}(z|h') \eta_{i,h'}(\mathbf{dt}_{-i})}{\mathbb{P}_{\sigma_i, \zeta_{-i}}^{\eta_{i,h'}}(z|h')} \\ &= \mathbb{P}_{\sigma_i, \zeta_{-i}}^{\eta_{i,h}}(h'|h) \int_{E_{-i}} \mathbb{P}_{\sigma_i, \zeta_{-i}}^{\eta_{i,h'}}(z|h') \eta_{i,h'}(\mathbf{dt}_{-i}) \\ &= \mu_i(h'|h) \mu_i(\{z\} \times E_{-i}|h'), \end{aligned}$$

where the second equality follows from condition (6). Finally, we need to show that  $\mu_i$  satisfies OAL. Fix  $h \in H$ ,  $a_i, b_i \in A_i(h)$ ,  $a_{-i} \in A_{-i}(h)$  and a measurable set  $E_{-i} \subseteq T_{-i}$ . We have

$$\begin{aligned} \mu_i(Z(h, a_{-i}) \times E_{-i}|h, a_i) &= \mu_i(Z(h, a) \times E_{-i}|h, a_i) = \sum_{z \in Z(h, a)} \mu_i(\{z\} \times E_{-i}|h, a_i) \\ &= \sum_{z \in Z(h, a)} \int_{E_{-i}} \mathbb{P}_{\sigma_i, \zeta_{-i}}^{\eta_{i,h}}(z|h, a_i) \eta_{i,h}(\mathbf{dt}_{-i}) \\ &= \sum_{z \in Z(h, a)} \int_{E_{-i}} \mathbb{P}_{\sigma_i, \zeta_{-i}}^{\eta_{i,h}}(z|h, a) \zeta_{-i}^{t_{-i}}(a_{-i}|h) \eta_{i,h}(\mathbf{dt}_{-i}) \\ &= \int_{E_{-i}} \left[ \sum_{z \in Z(h, a)} \mathbb{P}_{\sigma_i, \zeta_{-i}}^{\eta_{i,h}}(z|h, a) \right] \zeta_{-i}^{t_{-i}}(a_{-i}|h) \eta_{i,h}(\mathbf{dt}_{-i}) \\ &= \int_{E_{-i}} \zeta_{-i}^{t_{-i}}(a_{-i}|h) \eta_{i,h}(\mathbf{dt}_{-i}) \\ &= \int_{E_{-i}} \left[ \sum_{z \in Z(h, (b_i a_{-i}))} \mathbb{P}_{\sigma_i, \zeta_{-i}}^{\eta_{i,h}}(z|h, (b_i a_{-i})) \right] \zeta_{-i}^{t_{-i}}(a_{-i}|h) \eta_{i,h}(\mathbf{dt}_{-i}) \\ &= \sum_{z \in Z(h, (b_i a_{-i}))} \int_{E_{-i}} \mathbb{P}_{\sigma_i, \zeta_{-i}}^{\eta_{i,h}}(z|h, b_i) \eta_{i,h}(\mathbf{dt}_{-i}) \\ &= \sum_{z \in Z(h, (b_i a_{-i}))} \int_{E_{-i}} \mathbb{P}_{\sigma_i, \zeta_{-i}}^{\eta_{i,h}}(z|h, (b_i, a_{-i})) \zeta_{-i}^{t_{-i}}(a_{-i}|h) \eta_{i,h}(\mathbf{dt}_{-i}) \end{aligned}$$

$$\begin{aligned}
 &= \sum_{z \in Z(h, (b_i, a_{-i}))} \mu_i(\{z\} \times E_{-i} | h, b_i) = \mu_i(Z(h, (b_i, a_{-i})) \times E_{-i} | h, b_i) \\
 &= \mu_i(Z(h, a_{-i}) \times E_{-i} | h, b_i).
 \end{aligned}$$

This shows that OAI holds.  $\square$

**Proof of Lemma 2.** From Lemma 1 we know that  $\Delta^{H_i}(\Omega_{-i})$  is compact metrizable. Consider a sequence  $(\mu_i^n)_{n \in \mathbb{N}}$  of elements in  $\Delta_i^{H_i}(\Omega_{-i})$  converging to  $\mu_i \in \Delta^{H_i}(\Omega_{-i})$ . We need to show that  $\mu_i \in \Delta_i^{H_i}(\Omega_{-i})$ . For simplicity, we write  $\mu_{i, h_i}^n$  and  $\mu_{i, h_i}$  to denote the corresponding conditional probabilities at  $h_i \in H_i$ . Fix  $h \in H$ ,  $a_i, a'_i \in A_i(h)$  and  $a_{-i} \in A_{-i}(h)$ , and consider the following class of subsets of  $\mathcal{B}(T_{-i})$ :

$$\mathcal{D} = \left\{ E_{-i} \in \mathcal{B}(T_{-i}) : \mu_{i, (h, a_i)}(Z(h, (a_i, a_{-i})) \times E_{-i}) = \mu_{i, (h, a'_i)}(Z(h, (a'_i, a_{-i})) \times E_{-i}) \right\}.$$

We show that all the open subsets of  $T_{-i}$  belong to  $\mathcal{D}$  and that  $\mathcal{D}$  is a Dynkin class. Finally, by the Dynkin's lemma, we have  $\mathcal{B}(T_{-i}) = \mathcal{D}$ . Also, note that, for each  $b_i \in \{a_i, a'_i\}$ , the map

$$E_{-i} \mapsto \mu_{i, (h, b_i)}(Z(h, (b_i, a_{-i})) \times E_{-i})$$

is a finite measure on  $\mathcal{B}(T_{-i})$ . In what follows, we will consider integrals with respect to the measures just introduced. For the sake of simplicity, we denote such measures, for all  $E_{-i} \in \mathcal{B}(T_{-i})$ , as

$$\begin{aligned}
 \mu_n(E_{-i}) &= \mu_{i, (h, a_i)}^n(Z(h, (a_i, a_{-i})) \times E_{-i}), \\
 \mu(E_{-i}) &= \mu_{i, (h, a_i)}(Z(h, (a_i, a_{-i})) \times E_{-i}), \\
 \mu'_n(E_{-i}) &= \mu_{i, (h, a'_i)}^n(Z(h, (a'_i, a_{-i})) \times E_{-i}), \\
 \mu'(E_{-i}) &= \mu_{i, (h, a'_i)}(Z(h, (a'_i, a_{-i})) \times E_{-i}).
 \end{aligned}$$

We proceed by steps.

1. (For all  $n \in \mathbb{N}$ , and for all measurable functions  $f : T_{-i} \rightarrow \mathbb{R}$ ,  $\int_{T_{-i}} f d\mu_n = \int_{T_{-i}} f d\mu'_n$ )<sup>50</sup> Fix  $n \in \mathbb{N}$ . If  $f = \mathbf{I}_{E_{-i}}$  (i.e.,  $f$  is the indicator function on  $E_{-i}$ ) for some measurable set  $E_{-i} \in \mathcal{B}(T_{-i})$ , then the thesis is true since  $\mu_i^n \in \Delta_i^{H_i}(\Omega_{-i})$ . If  $f$  is a simple measurable function, then there exists a finite partition  $\{E_{-i}^1, \dots, E_{-i}^Q\}$  of  $T_{-i}$  and a collection of real numbers  $\{d^1, \dots, d^Q\}$  such that  $f = \sum_{q=1}^Q d^q \mathbf{I}_{E_{-i}^q}$ . Therefore,

$$\begin{aligned}
 \int_{T_{-i}} f d\mu_n &= \int_{T_{-i}} \left( \sum_{q=1}^Q d^q \mathbf{I}_{E_{-i}^q} \right) d\mu_n \\
 &= \sum_{q=1}^Q d^q \left( \int_{T_{-i}} \mathbf{I}_{E_{-i}^q} d\mu_n \right) \\
 &= \sum_{q=1}^Q d^q \left( \mu_{i, (h, a_i)}^n(Z(h, (a_i, a_{-i})) \times E_{-i}^q) \right) \\
 &= \sum_{q=1}^Q d^q \left( \mu_{i, (h, a'_i)}^n(Z(h, (a'_i, a_{-i})) \times E_{-i}^q) \right) \\
 &= \sum_{q=1}^Q d^q \left( \int_{T_{-i}} \mathbf{I}_{E_{-i}^q} d\mu'_n \right) \\
 &= \sum_{q=1}^Q d^q \left( \int_{T_{-i}} \mathbf{I}_{E_{-i}^q} d\mu'_n \right) = \int_{T_{-i}} (f) d\mu'_n.
 \end{aligned}$$

<sup>50</sup> Here, with an abuse of notation, we let  $\int_{T_{-i}} (\cdot) d\mu_{i, h}$  denote the integral of the marginal over  $T_{-i}$  conditional on  $h$ .



If  $f$  is an arbitrary measurable function, then there exists a sequence  $(f^m)_{m \in \mathbb{N}}$  of simple measurable functions such that  $f^m \uparrow f$ . Therefore,

$$\begin{aligned} \int_{T_{-i}} f d\mu_n &= \int_{T_{-i}} \left( \lim_m f^m \right) d\mu_n \\ &= \int_{T_{-i}} \lim_m f^m d\mu_n \\ &= \lim_m \int_{T_{-i}} f^m d\mu_n \\ &= \lim_m \int_{T_{-i}} f^m d\mu'_n \\ &= \int_{T_{-i}} \lim_m f^m d\mu'_n \\ &= \int_{T_{-i}} f d\mu'_n, \end{aligned}$$

where the third and fifth equalities follow from the Monotone Convergence Theorem. This proves the claim.

2. (Every open subset of  $T_{-i}$  is in  $\mathcal{D}$ .) Consider an open set  $E_{-i} \in \mathcal{D}$ . By Urysohn's lemma, there exists a sequence  $(f^m)_{m \in \mathbb{N}}$  of continuous real functions defined over  $T_{-i}$  such that  $f^m(t_{-i}) \uparrow \mathbf{1}_{E_{-i}}(t_{-i})$  for all  $t_{-i} \in T_{-i}$ . Then, we have

$$\begin{aligned} \mu(E_{-i}) &= \int_{T_{-i}} \mathbf{1}_{E_{-i}} d\mu \\ &= \int_{T_{-i}} \lim_m (f^m) d\mu \\ &= \lim_m \int_{T_{-i}} f^m d\mu \\ &= \lim_m \lim_n \int_{T_{-i}} f^m d\mu_n \\ &= \lim_m \lim_n \int_{T_{-i}} f^m d\mu'_n \\ &= \lim_m \int_{T_{-i}} f^m d\mu' \\ &= \int_{T_{-i}} \lim_m f^m d\mu' = \mu'(E_{-i}), \end{aligned}$$

where the third and seventh equalities follow from the Monotone Convergence Theorem, the fourth and sixth equality follow from the characterization of weak convergence of measures (see Portmanteau Theorem), and the fifth equality follows from point 1.

3. ( $\mathcal{D}$  is a Dynkin class.) It is immediate to show that  $T_{-i} \in \mathcal{D}$ . Let  $E_{-i}, E'_{-i} \in \mathcal{D}$  such that  $E_{-i} \subseteq E'_{-i}$ . With this,

$$\begin{aligned} \mu_{i,(h,a_i)}(Z(h, (a_i, a_{-i})) \times (E'_{-i} \setminus E_{-i})) &= \mu(E'_{-i} \setminus E_{-i}) \\ &= \mu(E'_{-i}) - \mu(E_{-i}) \\ &= \mu'(E'_{-i}) - \mu'(E_{-i}) \\ &= \mu'(E'_{-i} \setminus E_{-i}) \\ &= \mu_{i,(h,a'_i)}(Z(h, (a'_i, a_{-i})) \times (E'_{-i} \setminus E_{-i})), \end{aligned}$$

showing that  $(E'_{-i} \setminus E_{-i}) \in \mathcal{D}$ . Finally, consider a sequence  $(E^n_{-i})_{n \in \mathbb{N}}$  of pairwise disjoint measurable subsets of  $T_{-i}$  in  $\mathcal{D}$ . Then

$$\begin{aligned} \mu_{i,(h,a_i)} \left( Z(h, (a_i, a_{-i})) \times \left( \bigcup_n E^n_{-i} \right) \right) &= \sum_n \mu(E^n_{-i}) \\ &= \sum_n \mu'(E^n_{-i}) \\ &= \mu_{i,(h,a'_i)} \left( Z(h, (a'_i, a_{-i})) \times \left( \bigcup_n E^n_{-i} \right) \right), \end{aligned}$$

which shows that  $\bigcup E^n_{-i} \in \mathcal{D}$ . This finally shows that  $\mathcal{D}$  is a Dynkin class and  $\mathcal{D} = \mathcal{B}(T_{-i})$ .

Given that  $h \in H$ ,  $a_i, a'_i \in A_i(h)$ ,  $a_{-i} \in A_{-i}(h)$  were arbitrarily chosen,  $\mu_i$  satisfies OAI and belongs to  $\Delta_i^{H_i}(\Omega_{-i})$ , proving that the latter is closed, hence compact metrizable.  $\square$

**Proof of Lemma 3.** By Theorem 9 in Battigalli et al. (2017), there exists  $\nu_i \in \Delta_i^{H_i}(\Omega_{-i} \times X_{-i})$  that strongly believes  $(E^1_{-i}, \dots, E^n_{-i})$  and satisfies  $\text{marg}_{\Omega_{-i}} \nu_i = \mu_i$ . We need to show that  $\nu_i \in \Delta_i^{H_i}(Z \times T_{-i} \times X_{-i})$ . By inspection of the proof of Theorem 9 in Battigalli et al. (2017), we know that

$$\nu_i(B|h_i) = \mu_i^*(f^{-1}(B)|h_i)$$

for each measurable  $B \subseteq Z \times T_{-i} \times X_{-i}$ , where  $\mu_i^*(\cdot|h_i)$  is the completion of  $\mu_i(\cdot|h_i)$  and function  $f : Z \times T_{-i} \rightarrow Z \times T_{-i} \times X_{-i}$  is analytically measurable and defined as

$$f(z, t_{-i}) = (z, t_{-i}, q(z, t_{-i}))$$

for some analytically measurable  $q : Z \times T_{-i} \rightarrow Z \times T_{-i} \times X_{-i}$ .

Next, fix  $h \in H$ ,  $a_i, a'_i \in A_i(h)$ ,  $a_{-i} \in A_{-i}(h)$  and measurable  $E_{-i} \subseteq T_{-i} \times X_{-i}$ . We have

$$\mu_i(Z(h, (a_i, a_{-i})) \times E_{-i} \times B_{-i}|h, a_i) = \mu_i(Z(h, (a'_i, a_{-i})) \times E_{-i}|h, a'_i),$$

and

$$\begin{aligned} \nu_i(Z(h, (a_i, a_{-i})) \times E_{-i} \times B_{-i}|h, a_i) &= \mu_i^*(f^{-1}(Z(h, (a_i, a_{-i})) \times E_{-i} \times B_{-i})|h, a_i) \\ &= \mu_i(Z(h, (a_i, a_{-i})) \times E_{-i}|h, a_i) \\ &= \mu_i(Z(h, (a'_i, a_{-i})) \times E_{-i}|h, a'_i) \\ &= \mu_i^*(f^{-1}(Z(h, (a'_i, a_{-i})) \times E_{-i} \times B_{-i})|h, a'_i) \\ &= \nu_i(Z(h, (a'_i, a_{-i})) \times E_{-i} \times B_{-i}|h, a'_i), \end{aligned}$$

showing that also  $\nu_i$  satisfies OAI.  $\square$

**Proof of Proposition 2.** We need some preliminary definitions. Let  $Y_i^\infty$  be the set of coherent infinite hierarchies that do not necessarily satisfy OAI. Clearly, we have  $C_i^\infty \subseteq Y_i^\infty$ . From Proposition 1 of Battigalli and Siniscalchi (1999) we know that there exists a homeomorphism  $b_i : Y_i^\infty \rightarrow \Delta_i^{H_i}(\Omega_{-i}^\infty)$  such that

$$\mu_i^k = \text{marg}_{\Omega_{-i}^k} b_i(\mu_i^\infty) \tag{7}$$

for all  $k \in \mathbb{N}$ . We only need to check that  $b_i(C_i^\infty) = \Delta_i^{H_i}(\Omega_{-i}^\infty)$  and define  $g_i = b_i|_{C_i^\infty}$ . Consider  $\nu_i \in b_i(C_i^\infty) \subseteq b_i(Y_i^\infty) = \Delta_i^{H_i}(\Omega_{-i}^\infty)$ . Thus there exists  $\mu_i^\infty \in C_i^\infty$  such that  $b_i(\mu_i^\infty) = \nu_i$ . Fix any  $h \in H$ ,  $a_i, a'_i \in A_i(h)$  and  $a_{-i} \in A_{-i}(h)$ , and consider the families of subsets

$$\mathcal{D} = \left\{ E_{-i} \in \mathcal{B}(T_{-i}^\infty) : \nu_{i,(h,(a_i,a_{-i}))}(Z(h_i, (a_i, a_{-i})) \times E_{-i}) = \nu_{i,(h,(a'_i,a_{-i}))}(Z(h_i, (a'_i, a_{-i})) \times E_{-i}) \right\}$$

and the class of cylinder subsets of  $T_{-i}^\infty(h_i)$

$$\mathcal{C} = \left\{ E_{-i}^k \times T_{-i}^\infty \subseteq T_{-i}^\infty : k \in \mathbb{N}, E_{-i}^k \in \mathcal{B}(T_{-i}^k) \right\}.$$

Given that  $\mu_i^\infty \in C_i^\infty$  and (7) holds, it is easy to verify that  $\mathcal{C} \subseteq \mathcal{D}$ . With essentially the same steps used in the proof of Lemma 2 one can show that  $\mathcal{D}$  is a Dynkin class of subsets of  $T_{-i}^\infty$  and therefore, by Dynkin's Lemma and the well

known fact  $\mathcal{C}$  is a  $\pi$ -class, we have  $\mathcal{D} = \mathcal{B}(T_{-i}^\infty)$ . This finally shows that  $v_i \in \Delta_i^{H_i}(\Omega_{-i}^\infty)$ . Next, pick any  $v_i \in \Delta_i^{H_i}(\Omega_{-i}^\infty)$ . We want to show that there exists  $\mu_i^\infty \in C_i^\infty$  such that  $b_i(\mu_i^\infty) = v_i$ . Given that  $\Delta_i^{H_i}(\Omega_{-i}^\infty) \subseteq \Delta^{H_i}(\Omega_{-i}^\infty) = b_i(Y_i^\infty)$ , there exists  $\mu_i^\infty \in Y_i^\infty$  with  $b_i(\mu_i^\infty) = v_i$ . Finally, (7) and  $v_i \in \Delta_i^{H_i}(\Omega_{-i}^\infty)$  necessarily implies that each  $\mu_i^k$  satisfies OAI, showing that  $\mu_i^\infty \in C_i^\infty$ .  $\square$

**Proof of Lemma 5.** Define the correspondence

$$\Sigma_i^* : \begin{array}{l} \Omega_i^\infty \quad \Rightarrow \quad \Sigma_i, \\ (z, \theta_i, \mu_i^\infty) \mapsto \prod_{h \in H} \Delta(r_{i,h}(\theta_i, \mu_{i,k+1})). \end{array}$$

Note that  $\Sigma_i^*$  inherits all the properties of each  $r_{i,h}$  (in particular, it is upper hemicontinuous) and that

$$RP_i = \left\{ \omega_i^\infty \in \Omega_i^{\infty,*} : \omega_i^\infty \in (\hat{\sigma}_i)^{-1}(\Sigma_i^*(\omega_i^\infty)) \right\},$$

that is,  $RP_i$  coincides with the set of fixed points of the correspondence  $(\hat{\sigma}_i)^{-1} \circ \Sigma_i^*$ . By upper hemicontinuity of  $(\hat{\sigma}_i)^{-1} \circ \Sigma_i^*$  and Kakutani fixed point theorem, it follows that  $RP_i$  is non-empty and compact.  $\square$

**Proof of Lemma 6.** Let  $N$  denote the maximum length of the game. Note that  $N$  is well defined as  $Z$  is finite. Both  $MC_i$  and  $MC_i^*$  are obviously nonempty. Define the function  $q_i : \Omega_i^{\infty,*} \rightarrow \mathbb{R}^N$  as

$$q_i(z, \theta, \mu^\infty)_n = \begin{cases} g_i(\mu_i^\infty) \left( \left[ h_i^{n-1}(z), a_{i,n}(z) \right] \middle| h_i^{n-1}(z) \right), & \text{if } n \leq \ell(z), \\ c, & \text{otherwise,} \end{cases}$$

for all  $n \in \{1, \dots, N\}$ , where  $c \in \mathbb{R}$  is an arbitrary real number. In words, the function  $q_i(\cdot)$  takes value in those states of the world at which player  $i$  is fully coherent and gives back the list of probabilities with which he planned to play the actions implied by  $z$ . With this, we can write

$$\begin{aligned} MC_i &= \{ (z, \theta, \mu^\infty) \in \Omega_i^{\infty,*} : q_i(z, \theta, \mu^\infty) > \mathbf{0} \}, \\ sMC_i &= \{ (z, \theta, \mu^\infty) \in \Omega_i^{\infty,*} : q_i(z, \theta, \mu^\infty) = \mathbf{1} \}, \end{aligned}$$

where  $\mathbf{0} = (0, \dots, 0)$ ,  $\mathbf{1} = (1, \dots, 1) \in \mathbb{R}^N$ . We thus need to show that  $q$  is measurable with respect to the Borel sigma-algebra over  $\Omega_i^{\infty,*}$ . In particular, it is sufficient (and necessary) that  $q_i(\cdot)_n : \Omega_i^{\infty,*} \rightarrow \mathbb{R}$  is measurable for each  $n \in \{1, \dots, N\}$ . Recalling that (1)  $g_i$  is a homeomorphism, (2)  $h_i^{n-1}(\cdot)$  and  $a_{i,n}(\cdot)$  are functions between finite spaces, and (3) the Borel sigma-algebra of the set of probability measures (endowed with the topology of weak convergence) on a compact metrizable space is generated by all the bounded continuous functionals over that space, we can conclude that each  $q_i(\cdot)_n$  is continuous, hence measurable. Continuity of each  $q_i(\cdot)_n$  also implies that  $MC_i^*$  is closed, hence compact.  $\square$

**Proof of Proposition 3.** We prove the result by induction. Clearly,

$$\mathbf{P}_i(1) \subseteq \Omega_i^k = \mathbf{P}_i(0), \mathbf{P}_{-i}(1) \subseteq \Omega_{-i}^k = \mathbf{P}_{-i}(0) \text{ and } \mathbf{P}(1) \subseteq \Omega^k = \mathbf{P}(0).$$

Moreover, note that  $\mathbf{P}_i(1)$  can be written as the intersection of the following two sets:

$$\left\{ (z, \theta_i, \mu_i^k) \in \Omega_i^k : \left( \mu_i^k \in C_i^k \right) \wedge \left( \forall h \in H, h < z \implies \hat{\sigma}_i(\mu_{i,1})(a_{i,h}(z) | h) > 0 \right) \right\}$$

and

$$\text{proj}_{\Omega_i^k} \left\{ (z, \theta_i, \mu_i^{k+1}) \in \Omega_i^{k+1} : \left( \mu_i^{k+1} \in C_i^{k+1} \right) \wedge \left( \forall h \in H, \hat{\sigma}_i(\mu_{i,1})(r_{i,h}(\theta_i, \mu_{i,k+1}) | h) = 1 \right) \right\}.$$

On the one hand, through the same steps used in the proof of Lemma 6, one can show that the former set is measurable. On the other hand, the latter set is the image through a continuous function (i.e., the projection) of a compact set, hence measurable. With this,  $\mathbf{P}_i(1)$  is measurable as well. Next, assume that (i) – (ii) hold for every  $k \in \{0, 1, \dots, n\}$ . Let  $(z, \theta_i, \mu_i^k) \in \mathbf{P}_i(n+1)$ . It follows that there exists  $\mu_{i,k+1} \in M_{i,k+1}$  such that  $(z, \theta_i, \mu_i^k)$  and  $\mu_{i,k+1}$  satisfy Coherence, RP, MC and Strong belief for each  $m \in \{1, \dots, n\}$ . Therefore,  $(z, \theta_i, \mu_i^k) \in \mathbf{P}_i(n)$ . We can similarly show that

$$\mathbf{P}_{-i}(n+1) \subseteq \mathbf{P}_{-i}(n) \text{ and } \mathbf{P}(n+1) \subseteq \mathbf{P}(n).$$

For (i), measurability of  $\mathbf{P}_i(n+1)$  follows from the fact that  $\mathbf{P}_i(n)$  is measurable and the measurability property of strong belief (see Lemma 6).  $\square$

**Proof of Theorem 1.** We prove the thesis by induction on  $n$ .

**(Basis step,  $n = 1$ )** Fix  $i \in I$ . Pick any  $(z, \theta_i, \mu_i^k) \in \mathbf{P}_i(1)$ . Then, there exists  $\mu_{i,k+1} \in \Delta_i^{H_i}(\Omega_{-i}^k)$  such that  $(z, \theta_i, \mu_i^k, \mu_{i,k+1})$  satisfies coherence, RP and MC. It is not hard to verify that  $(z, \theta_i, \bar{\mu}_i^\infty) \in R_i(1)$  for all  $\bar{\mu}_i^\infty \in C_i^\infty$  such that  $\bar{\mu}_i^{k+1} = (\mu_i^k, \mu_{i,k+1})$ , and therefore  $(z, \theta_i, \mu_i^k) \in \text{proj}_{\Omega_i^k} R_i(1)$ . Conversely, let  $(z, \theta_i, \mu_i^k) \in \text{proj}_{\Omega_i^k} R_i(1)$ . Then, there exists  $\bar{\mu}_i^\infty \in C_i^\infty$  such that  $(z, \theta_i, \bar{\mu}_i^\infty) \in R_i(1)$  and  $\bar{\mu}_i^k = \mu_i^k$ . One can check that  $\bar{\mu}_{i,k+1} \in \Delta_i^{H_i}(\Omega_{-i}^k)$  is such that  $(z, \theta_i, \mu_i^k, \bar{\mu}_{i,k+1})$  satisfies coherence, RP and MC, that is,  $(z, \theta_i, \mu_i^k) \in \mathbf{P}_i(1)$ . Since  $i$  was arbitrarily chosen, it follows that  $\mathbf{P}_i(1) = \text{proj}_{\Omega_i^k} R_i(1)$  for every  $i \in I$ .

**(Inductive step)** Assume that  $\mathbf{P}_{-i}(m) = \text{proj}_{\Omega_i^k} R_{-i}(m)$  for every  $i \in I$  and  $m \in \{1, \dots, n\}$ . First, we show that the inductive hypothesis implies that  $\mathbf{P}_i(m) = \text{proj}_{\Omega_i^k} R_{-i}(m)$  for every  $i \in I$  and  $m \in \{1, \dots, n\}$ . For the sake of simplicity, we write  $T_{-i,j}^k$  instead of  $\prod_{t \in I \setminus \{i,j\}} T_t^k$ , for every  $i, j \in I$  and  $k \in \mathbb{N} \cup \{\infty\}$ . Fix  $i \in I$  and  $m \in \{0, \dots, n\}$ . Pick any  $(z, \theta_{-i}, \mu_{-i}^k) \in \mathbf{P}_{-i}(m) = \bigcap_{j \in I \setminus \{i\}} (\mathbf{P}_j(m) \times T_{-i,j}^k)$ . Then  $(z, \theta_j, \mu_j^k) \in \mathbf{P}_j(m)$  for every  $j \in I \setminus \{i\}$ . By the inductive hypothesis, for every  $j \in I \setminus \{i\}$ , there exists  $\bar{\mu}_j^\infty \in M_j^\infty$  such that  $\bar{\mu}_j^k = \mu_j^k$  and  $(z, \theta_j, \bar{\mu}_j^\infty) \in R_j(m)$ . Moreover, by definition it holds that  $R_{-i}(m) = \bigcap_{j \in I \setminus \{i\}} (R_j(m) \times T_{-i,j}^\infty)$ , hence, for every  $j \in I \setminus \{i\}$ , we have  $(z, \theta_{-i}, \bar{\mu}_{-i}^\infty) \in R_j(m) \times T_{-i,j}^\infty$ , where  $\bar{\mu}_{-i}^\infty = (\bar{\mu}_j^\infty)_{j \in I \setminus \{i\}}$ . Therefore,  $(z, \theta_{-i}, \bar{\mu}_{-i}^\infty) \in R_{-i}(m)$ , proving that  $(z, \theta_{-i}, \mu_{-i}^k) \in \text{proj}_{\Omega_i^k} R_{-i}(m)$ . Conversely, assume that  $(z, \theta_{-i}, \mu_{-i}^k) \in \text{proj}_{\Omega_i^k} R_{-i}(m)$ . It follows that there exists  $\bar{\mu}_{-i}^\infty$  such that  $\bar{\mu}_{-i}^k = \mu_{-i}^k$  and, for every  $j \in I \setminus \{i\}$ ,  $(z, \theta_j, \bar{\mu}_j^\infty) \in R_j(m)$ . By the inductive hypothesis, for every  $j \in I \setminus \{i\}$ , we have  $(z, \theta_j, \mu_j^k) \in \mathbf{P}_j(m)$  and, as a consequence,  $(z, \theta_{-i}, \mu_{-i}^k) \in \mathbf{P}_j(m) \times T_{-i,j}^k$ . With this,  $(z, \theta_{-i}, \mu_{-i}^k) \in \mathbf{P}_{-i}(m)$ . Since  $i$  and  $m$  were arbitrarily chosen, the claim holds.

Next, we show that  $\mathbf{P}_i(n+1) = \text{proj}_{\Omega_i^k} R_i(n+1)$  for every  $i \in I$ . Fix  $i \in I$ , and assume first that  $(z, \theta_i, \mu_i^k) \in \mathbf{P}_i(n+1)$ . Then, there exists  $\mu_{i,k+1} \in \Delta_i^{H_i}(\Omega_{-i}^k)$  such that  $(z, \theta_i, \mu_i^k, \mu_{i,k+1})$  satisfies coherence, RP and MC. Moreover,  $\mu_{i,k+1}$  strongly believes the decreasing chain  $(\mathbf{P}_{-i}(m))_{m=1}^n$  of events in  $\Omega_{-i}^k$ . By the previous claim, we have that  $\mu_{i,k+1}$  strongly believes  $(\text{proj}_{\Omega_i^k} R_{-i}(m))_{m=1}^n$ . Therefore, by Lemma 3, there exists  $\nu_i \in \Delta_i^{H_i}(\Omega_{-i}^\infty)$  that strongly believes the decreasing chain  $(R_{-i}(m))_{m=1}^n$  of events in  $\Omega_{-i}^\infty$  and such that  $\text{marg}_{\Omega_{-i}^k} \nu_i = \mu_{i,k+1}$ . Next, let  $\bar{\mu}_i^\infty \in C_i^\infty$  be defined as  $g_i^{-1}(\nu_i)$ . We claim that  $(z, \theta_i, \bar{\mu}_i^\infty) \in R_i(n+1)$  and  $\bar{\mu}_i^k = \mu_i^k$ . The second part is immediate since, by Proposition 2, for all  $q \leq k$ ,

$$\bar{\mu}_{i,q} = \text{marg}_{\Omega_{-i}^{q-1}} g_i(\bar{\mu}_i^\infty) = \text{marg}_{\Omega_{-i}^{q-1}} g_i(g_i^{-1}(\nu_i)) = \text{marg}_{\Omega_{-i}^{q-1}} \nu_i = \text{marg}_{\Omega_{-i}^{q-1}} \mu_{i,k+1} = \mu_{i,q},$$

where the last equality follows from the fact that  $(\mu_i^k, \mu_{i,k+1}) \in C_i^{k+1}$  by hypothesis. As for the first part of the claim, note that

$$R_i(n+1) = R_i(n) \cap \text{SB}_i(R_{-i}(n)) = R_i \cap \bigcap_{m=1}^n \text{SB}_i(R_{-i}(m)).$$

Therefore, it is enough to show that  $(z, \theta_i, \bar{\mu}_i^\infty) \in R_i$  and  $(z, \theta_i, \bar{\mu}_i^\infty) \in \text{SB}_i(R_{-i}(m))$  for every  $m \in \{1, \dots, n\}$ . The fact that  $(z, \theta_i, \bar{\mu}_i^\infty) \in R_i$  is trivial since  $\bar{\mu}_i^\infty \in C_i^\infty$  and  $\bar{\mu}_i^k = \mu_i^k$  satisfies rational planning and material consistency. Since  $g_i(\bar{\mu}_i^\infty) = \nu_i$  strongly believes  $R_{-i}(m)$  for every  $m \in \{1, \dots, n\}$ , it follows that  $(z, \theta_i, \bar{\mu}_i^\infty) \in \text{SB}_i(R_{-i}(m))$  for every  $m \in \{1, \dots, n\}$ . With this, we proved that  $(z, \theta_i, \bar{\mu}_i^\infty) \in R_i(n+1)$  and therefore that  $(z, \theta_i, \mu_i^k) = (z, \theta_i, \bar{\mu}_i^k) \in \text{proj}_{\Omega_i^k} R_i(n+1)$ . Conversely, assume that  $(z, \theta_i, \mu_i^k) \in \text{proj}_{\Omega_i^k} R_i(n+1)$ . Thus there exists  $\bar{\mu}_i^\infty \in M_i^\infty$  such that  $(z, \theta_i, \bar{\mu}_i^\infty) \in R_i(n+1)$  and  $\bar{\mu}_i^k = \mu_i^k$ . Now, consider  $\bar{\mu}_{i,k+1} \in \Delta_i^{H_i}(\Omega_{-i}^k)$ . It is clear that  $(\mu_i^k, \bar{\mu}_{i,k+1}) \in C_i^{k+1}$  and that  $(z, \theta_i, \mu_i^k, \bar{\mu}_{i,k+1})$  satisfies RP and MC since  $(z, \theta_i, \bar{\mu}_i^\infty) \in R_i(n+1) \subseteq R_i$ . We still need to show that  $\bar{\mu}_{i,k+1}$  strongly believes the chain  $(P_{-i}(m))_{m=1}^n$ . By the first claim, this is equivalent to showing that  $\bar{\mu}_{i,k+1}$  strongly believes the chain  $(\text{proj}_{\Omega_i^k} R_{-i}(m))_{m=1}^n$ . Pick any  $h \in H$  and  $m \in \{1, \dots, n\}$  such that  $\Omega_{-i}^k(h) \cap \text{proj}_{\Omega_i^k} R_{-i}(m) \neq \emptyset$ . Then, there exists  $z \in Z(h)$  and  $t_{-i}^\infty \in T_{-i}^\infty$  such that  $(z, t_{-i}^\infty) \in R_{-i}(m)$ . Then, we have  $\Omega_{-i}^\infty(h) \cap R_{-i}(m) \neq \emptyset$  which, by hypothesis, implies  $g_i(\bar{\mu}_i^\infty)(R_{-i}(m) | h) = 1$ . By Proposition 2,

$$\begin{aligned} \bar{\mu}_{i,k+1}(\text{proj}_{\Omega_i^k} R_{-i}(m) | h) &= \text{marg}_{\Omega_{-i}^k} g_i(\bar{\mu}_i^\infty) (\text{proj}_{\Omega_i^k} R_{-i}(m) | h) \\ &= g_i(\bar{\mu}_i^\infty) \left( \left( \text{proj}_{\Omega_{-i}^k} \right)^{-1} \circ \left( \text{proj}_{\Omega_{-i}^k} \right) (R_{-i}(m)) | h \right) \\ &= g_i(\bar{\mu}_i^\infty) (R_{-i}(m) | h) = 1. \end{aligned}$$

Given that  $h$  and  $m$  were arbitrarily chosen, this implies that  $\bar{\mu}_{i,k+1}$  strongly believes the chain  $\left(\text{proj}_{\Omega_{-i}^k} R_{-i}(m)\right)_{m=1}^n$ , i.e.,  $(P_{-i}(m))_{m=1}^n$ , showing that  $(z, \theta_i, \mu_i^k) \in \mathbf{P}_i(n+1)$ . Since  $i$  was arbitrarily chosen,  $\mathbf{P}_i(n+1) = \text{proj}_{\Omega_i^k} R_i(n+1)$  for every  $i \in I$ . Since the result holds for all  $n$ ,

$$\mathbf{P}_i(\infty) = \bigcap_{n \in \mathbb{N}} \mathbf{P}_i(n) = \bigcap_{n \in \mathbb{N}} \text{proj}_{\Omega_i^k} R_i(n).$$

Finally, we need to show that  $\mathbf{P}(n) = \text{proj}_{\Omega^k} R(n)$  for every  $n \in \mathbb{N}$ . Fix  $n \in \mathbb{N}$ , and let  $(z, t^k) \in \mathbf{P}(n)$ . Then  $(z, t_i^k) \in \mathbf{P}_i(n)$  for every  $i \in I$ . Hence there exists  $\bar{t}_i^\infty \in T_i^\infty$  such that  $(z, \bar{t}_i^\infty) \in R_i(n)$  and  $\bar{t}_i^k = t_i^k$ . With this,  $(z, \bar{t}^\infty) \in R_i(n) \times T_{-i}^\infty$  for every  $i \in I$ , which implies  $(z, \bar{t}^\infty) \in \bigcap_{i \in I} (R_i(n) \times T_{-i}^\infty) = R(n)$ . This shows that  $(z, t^k) \in \text{proj}_{\Omega^k} R(n)$ . The proof of the converse is almost identical to the proof of the first claim and left to the reader.  $\square$

**Proof of Proposition 5.** By repeating the same steps of the proof of Theorem 1, we can show that, for every  $i \in I$ ,

$$\forall n \in \mathbb{N}, \mathbf{P}_i^*(n) = \text{proj}_{\Omega_i^k} R_i^*(n) \quad \text{and} \quad \mathbf{P}^*(n) = \text{proj}_{\Omega^k} R^*(n).$$

This, together with Proposition 3, implies that

$$\begin{aligned} \mathbf{P}^*(\infty) &= \bigcap_{n \in \mathbb{N}} \mathbf{P}^*(n) \\ &= \bigcap_{n \in \mathbb{N}} \text{proj}_{\Omega^k} R^*(n) \\ &\supseteq \text{proj}_{\Omega^k} \bigcap_{n \in \mathbb{N}} R^*(n) \\ &= \text{proj}_{\Omega^k} R^*(\infty). \end{aligned}$$

Conversely, let  $(z, \theta, \mu^k) \in \mathbf{P}^*(\infty) = \bigcap_{n \in \mathbb{N}} \text{proj}_{\Omega^k} R^*(n)$ , so that, for every  $n \in \mathbb{N}$ , there exists  $\bar{\mu}^\infty(n) \in M^\infty$  such that  $(z, \theta, \bar{\mu}^\infty(n)) \in R^*(n)$  and  $\bar{\mu}^k(n) = \mu$ . This implies that, for every  $n \in \mathbb{N}$ , the section  $(R^*(n))_{(z, \theta, \mu^k)}$  is nonempty. In particular,  $\left((R^*(n))_{(z, \theta, \mu^k)}\right)_{n \in \mathbb{N}}$  is a decreasing sequence of nonempty compact sets and, by the finite intersection property,  $\bigcap_{n \in \mathbb{N}} (R^*(n))_{(z, \theta, \mu^k)} \neq \emptyset$ . With this, pick any  $(\hat{\mu}^\ell)_{\ell \geq k+1} \in \bigcap_{n \in \mathbb{N}} (R^*(n))_{(z, \theta, \mu^k)}$ ; then we have, by construction,

$$\left(z, \theta, \mu^k, (\hat{\mu}^\ell)_{\ell \geq k+1}\right) \in \bigcap_{n \in \mathbb{N}} (R^*(n)) = R^*(\infty),$$

and so  $(z, \theta, \mu^k) \in \text{proj}_{\Omega^k} R^*(\infty)$ .  $\square$

## References

- Aliprantis, C., Border, K., 2006. *Infinite Dimensional Analysis*. Springer-Verlag, Berlin.
- Battigalli, P., Corrao, R., 2019. Iterated Dominance in Psychological Games. Bocconi University. Typescript.
- Battigalli, P., Dufwenberg, M., 2009. Dynamic psychological games. *J. Econ. Theory* 144, 1–35.
- Battigalli, P., Dufwenberg, M., 2019. Psychological Game Theory. IGIER Working Paper 646. Bocconi University.
- Battigalli, P., Siniscalchi, M., 1999. Hierarchies of conditional beliefs and interactive epistemology in dynamic games. *J. Econ. Theory* 88, 188–230.
- Battigalli, P., Siniscalchi, M., 2002. Strong belief and forward induction reasoning. *J. Econ. Theory* 106, 356–391.
- Battigalli, P., Tebaldi, P., 2019. Interactive epistemology in simple dynamic games with a continuum of strategies. *Econ. Theory* 68, 737–763.
- Battigalli, P., Beneduci, G., Tebaldi, P., 2017. Interactive Epistemology in Simple Dynamic Games with a Continuum of Strategies. IGIER Working Paper 602. Bocconi University.
- Battigalli, P., Corrao, R., Dufwenberg, M., 2019a. Incorporating belief-dependent motivation in games. *J. Econ. Behav. Organ.* 167, 185–218.
- Battigalli, P., Dufwenberg, M., Smith, A., 2019b. Frustration, aggression and anger in leader-follower games. *Games Econ. Behav.* 117, 15–39.
- Battigalli, P., Di Tillio, A., Samet, D., 2013. Strategies and interactive beliefs in dynamic games. In: Acemoglu, D., Arellano, M., Dekel, E. (Eds.), *Advances in Economics and Econometrics*. Cambridge University Press, Cambridge, UK, pp. 391–422.
- Brandenburger, A., Dekel, E., 1993. Hierarchies of beliefs and common knowledge. *J. Econ. Theory* 59, 189–198.
- Dekel, E., Siniscalchi, M., 2015. Epistemic game theory. In: Young, P., Zamir, S. (Eds.), *Handbook of Game Theory with Economic Applications*, vol. 4. North-Holland, Amsterdam, pp. 619–702.
- Geanakoplos, J., Pearce, D., Stacchetti, E., 1989. Psychological games and sequential rationality. *Games Econ. Behav.* 1, 60–79.
- Gul, F., Pesendorfer, W., 2016. Interdependent preference models as a theory of intentions. *J. Econ. Theory* 165, 179–208.
- Jagau, S., Perea, A., 2017. Common belief in rationality in psychological games. Epicenter Working Paper 10.
- Jagau, S., Perea, A., 2018. Expectation-based psychological games and psychological expected utility. Typescript.

- Kuhn, H.W., 1953. Extensive games and the problem of information. In: Kuhn, H.W., Tucker, A.W. (Eds.), *Contributions to the Theory of Games II*. Princeton University Press, Princeton, pp. 193–216.
- Levine, D.K., 1998. Modeling altruism and spitefulness in experiments. *Rev. Econ. Dyn.* 1, 593–622.
- Mourmans, N., 2019. Reasoning in Psychological Games: When is iterated Elimination of Choices Enough? Epicenter Working Paper 20. Maastricht University.
- Perea, A., 2012. *Epistemic Game Theory: Reasoning and Choice*. CUP Press.
- Shimoji, M., Watson, J., 1998. Conditional dominance, rationalizability, and game forms. *J. Econ. Theory* 83, 161–195.