

---

# A Logical Characterization of Iterated Admissibility

---

Joseph Y. Halpern and Rafael Pass

Department of Computer Science

Cornell University

Ithaca, NY, 14853, U.S.A.

e-mail: halpern@cs.cornell.edu, rafael@cs.cornell.edu

## Abstract

Brandenburger, Friedenberg, and Keisler provide an epistemic characterization of iterated admissibility (i.e., iterated deletion of weakly dominated strategies) where uncertainty is represented using LPSs (lexicographic probability sequences). Their characterization holds in a rich structure called a *complete* structure, where all types are possible. Here, a logical characterization of iterated admissibility is given that involves only standard probability and holds in all structures, not just complete structures. Roughly speaking, our characterization shows that iterated admissibility captures the intuition that “all the agent knows” is that agents satisfy the appropriate rationality assumptions.

## 1 Introduction

*Admissibility* is an old criterion in decision making. A strategy for player  $i$  is admissible if it is a best response to some belief of player  $i$  that puts positive probability on all the strategy profiles for the other players. Part of the interest in admissibility comes from the observation (due to Pearce [1984]) that a strategy  $\sigma$  for player  $i$  is admissible iff it is not weakly dominated; that is, there is no strategy  $\sigma'$  for player  $i$  that gives  $i$  at least as high a payoff as  $\sigma$  no matter what strategy the other players are using, and sometimes gives  $i$  a higher payoff.

It seems natural to ignore strategies that are not admissible. But there is a conceptual problem when it comes to dealing with *iterated* admissibility (i.e., iterated deletion of weakly dominated strategies). As

Mas-Colell, Whinston, and Green [1995, p. 240] put it:

[T]he argument for deletion of a weakly dominated strategy for player  $i$  is that he contemplates the possibility that every strategy combination of his rivals occurs with positive probability. However, this hypothesis clashes with the logic of iterated deletion, which assumes, precisely, that eliminated strategies are not expected to occur.

Brandenburger, Friedenberg, and Keisler [2008] (BFK from now on) resolve this paradox in the context of iterated deletion of weakly dominated strategies by assuming that strategies are not really eliminated. Rather, they assumed that strategies that are weakly dominated occur with infinitesimal (but nonzero) probability. (Formally, this is captured by using an LPS—*lexicographically ordered probability sequence*.) They define a notion of belief (which they call *assumption*) appropriate for their setting, and show that strategies that survive  $k$  rounds of iterated deletion are ones that are played in states where there is  $k$ th-order mutual belief in rationality; that is, everyone assumes that everyone assumes ... ( $k - 1$  times) that everyone is rational. However, they prove only that their characterization of iterated admissibility holds in particularly rich structures called *complete* structures (defined formally in Section 5), where all types are possible.

Here, we provide an alternate logical characterization of iterated admissibility. The characterization has the advantage that it holds in all structures, not just complete structures, and assumes that agents represent their uncertainty using standard probability measures, rather than LPSs or nonstandard probability measures

(as is done in a characterization of Rajan [1998]). Moreover, while complete structures must be uncountable, we show that our formula characterizing iterated admissibility is satisfiable in a structure with finitely many states.

Roughly speaking, instead of assuming only that agents know (or assume) that all other agents satisfy appropriate levels of rationality, we assume that “all the agents know” is that the other agents satisfy the appropriate rationality assumptions. We are using the phrase “all agent  $i$  knows” here in essentially the same sense that it is used by Levesque [1990] and Halpern and Lakemeyer [2001]. We formalize this notion by requiring that the agent ascribes positive probability to all formulas of some language  $\mathcal{L}$  that are consistent with his rationality assumptions. (This admittedly fuzzy description is made precise in Section 4.) As we show, when the language  $\mathcal{L}$  is sufficiently rich, our logical formula characterizes iterated admissibility. Slightly more precisely, let  $RAT_i$  be true iff player  $i$  is playing a best response to his belief. Define  $RAT_i^{k+1}$  to be true iff  $RAT_i$  holds, player  $i$  is playing a strategy constant with  $RAT_i^k$ , and for all  $j \neq i$ ,  $i$  knows that  $RAT_j^k$  holds, and that is “all that agent  $i$  knows” about players  $j \neq i$ . That is,  $RAT_i^{k+1}$  holds (i.e., player  $i$  is  $k+1$ -level rational) iff player  $i$  is playing a best response to his beliefs, using a strategy consistent with  $k$ -level rationality, and the only thing it knows about the other players is that they are  $k$ -level rational. As we show, for natural choices of languages  $\mathcal{L}$  a strategy  $\sigma$  for player  $i$  survives  $k$  levels of iterated deletion of weakly dominated strategies iff there is a structure and a state where  $\sigma$  is played by player  $i$  and the formula  $RAT_i^k$  holds.

For this result to hold,  $\mathcal{L}$  must be reasonably expressive; in particular, it must be possible to express in  $\mathcal{L}$  a player’s beliefs about the strategies that other players are playing. Interestingly, for less expressive languages  $\mathcal{L}$ ,  $RAT_i^k$  instead characterizes iterated deletion of *strongly* dominated strategies. For instance, if  $\mathcal{L}$  is empty, then “all agent  $i$  knows is  $\varphi$ ” is equivalent to “agent  $i$  knows  $\varphi$ ”; it then easily follows (given standard assumptions about knowledge) that  $RAT_i^{k+1}$  is equivalent to the statement that player  $i$  is rational and knows that everyone is rational and knows that everyone knows that everyone knows  $\dots$  ( $k-1$  times) that everyone is rational. Put another way, player  $i$  is  $k+1$  level rational iff player  $i$  is playing a best response to his beliefs, and he knows that all other players are  $k$ -

level rational.<sup>1</sup> Thus, essentially the same logical formula characterizes both iterated removal of strongly and weakly dominated strategies; in a sense, the only difference between these notions is the expressiveness of the language used by the players to reason about each other.

While we can take the structure where our characterizations hold to be countable, perhaps the most natural structure to consider is the *canonical* structure, which has a state corresponding to every satisfiable collection of formulas. The canonical structure is uncountable. We can show that the canonical structure is complete in the sense of BFK. Moreover, under a technical assumption, every complete structure is essentially canonical (i.e., it has a state corresponding to every satisfiable collection of formulas). This sequence of results allows us to connect iterated admissibility, complete structures, canonical structures, and the notion of “all I know”.

The rest of the paper is organized as follows. In Section 2 we introduce the formal model (which is essentially the standard Kripke model for probability, adapted to the game-theoretic setting), and provide some initial characterizations of rationalizability and iterated deletion of strongly dominated strategies. Section 3 and 4 contain our main characterization of iterated admissibility using “all I know”. In Section 3 we define “all I know” using a simple language; Section 4 considers the effect of using more expressive languages. Finally, in Section 5 we compare our results to those of BFK, and conclude with a discussion in Section 6.

## 2 Probability Structures, Rationalizability, and Admissibility

We consider normal-form games with  $n$  players. Given a (normal-form)  $n$ -player game  $\Gamma$ , let  $\Sigma_i(\Gamma)$  denote the strategies of player  $i$  in  $\Gamma$ . We omit the parenthetical  $\Gamma$  when it is clear from context or irrelevant. Let  $\vec{\Sigma} = \Sigma_1 \times \dots \times \Sigma_n$ .

Let  $\mathcal{L}^1$  be the language where we start with *true* and the special primitive proposition  $RAT_i$  and close off under modal operators  $B_i$  and  $\langle B_i \rangle$ , for  $i = 1, \dots, n$ , conjunction, and negation. We think of  $B_i\varphi$  as saying

<sup>1</sup>It follows from results of Tan and Werlang [Tan and Werlang 1988] that this formula characterizes rationalizability, and hence, by results of Pearce [1984], that it also characterizes iterated deletion of strongly dominated strategies.

that  $\varphi$  holds with probability 1, and  $\langle B_i \rangle \varphi$  as saying that  $\varphi$  holds with positive probability. As we shall see,  $\langle B_i \rangle$  is definable as  $\neg B_i \neg$  if we make the appropriate measurability assumptions.

To reason about the game  $\Gamma$ , we consider a class of probability structures corresponding to  $\Gamma$ . A *probability structure  $M$  appropriate for  $\Gamma$*  is a tuple  $(\Omega, \mathbf{s}, \mathcal{F}, \mathcal{PR}_1, \dots, \mathcal{PR}_n)$ , where  $\Omega$  is a set of states;  $\mathbf{s}$  associates with each state  $\omega \in \Omega$  a pure strategy profile  $\mathbf{s}(\omega)$  in the game  $\Gamma$ ;  $\mathcal{F}$  is a  $\sigma$ -algebra over  $\Omega$ ; and, for each player  $i$ ,  $\mathcal{PR}_i$  associates with each state  $\omega$  a probability distribution  $\mathcal{PR}_i(\omega)$  on  $(\Omega, \mathcal{F})$ . Intuitively,  $\mathbf{s}(\omega)$  is the strategy profile used at state  $\omega$  and  $\mathcal{PR}_i(\omega)$  is player  $i$ 's probability distribution at state  $\omega$ . As is standard, we require that each player knows his strategy and his beliefs. Formally, we require that (1) for each strategy  $\sigma_i$  for player  $i$ ,  $\llbracket \sigma_i \rrbracket_M = \{\omega : \mathbf{s}_i(\omega) = \sigma_i\} \in \mathcal{F}$ , where  $\mathbf{s}_i(\omega)$  denotes player  $i$ 's strategy in the strategy profile  $\mathbf{s}(\omega)$ ; (2)  $\mathcal{PR}_i(\omega)(\llbracket \mathbf{s}_i(\omega) \rrbracket_M) = 1$ ; (3) for each probability measure  $\pi$  on  $(\Omega, \mathcal{F})$ , and player  $i$ ,  $\llbracket \pi, i \rrbracket_M = \{\omega : \Pi_i(\omega) = \pi\} \in \mathcal{F}$ ; and (4)  $\mathcal{PR}_i(\omega)(\llbracket \mathcal{PR}_i(\omega), i \rrbracket_M) = 1$ .

The semantics is given as follows:

- $(M, \omega) \models \text{true}$  (so *true* is vacuously true).
- $(M, \omega) \models RAT_i$  if  $\mathbf{s}_i(\omega)$  is a best response, given player  $i$ 's beliefs on the strategies of other players induced by  $\mathcal{PR}_i(\omega)$ . (Because we restrict to appropriate structures, a player's expected utility at a state  $\omega$  is well defined, so we can talk about best responses.)
- $(M, \omega) \models \neg \varphi$  if  $(M, \omega) \not\models \varphi$ .
- $(M, \omega) \models \varphi \wedge \varphi'$  iff  $(M, \omega) \models \varphi$  and  $(M, \omega) \models \varphi'$ .
- $(M, \omega) \models B_i \varphi$  if there exists a set  $F \in \mathcal{F}_i$  such that  $F \subseteq \llbracket \varphi \rrbracket_M$  and  $\mathcal{PR}_i(\omega)(F) = 1$ , where  $\llbracket \varphi \rrbracket_M = \{\omega : (M, \omega) \models \varphi\}$ .
- $(M, \omega) \models \langle B_i \rangle \varphi$  if there exists a set  $F \in \mathcal{F}_i$  such that  $F \subseteq \llbracket \varphi \rrbracket_M$  and  $\mathcal{PR}_i(\omega)(F) > 0$ .

Given a language (set of formulas)  $\mathcal{L}$ ,  $M$  is  $\mathcal{L}$ -*measurable* if  $M$  is appropriate (for some game  $\Gamma$ ) and  $\llbracket \varphi \rrbracket_M \in \mathcal{F}$  for all formulas  $\varphi \in \mathcal{L}$ . It is easy to check that in an  $\mathcal{L}^1$ -measurable structure,  $\langle B_i \rangle \varphi$  is equivalent to  $\neg B_i \neg \varphi$ .

To put our results on iterated admissibility into context, we first consider rationalizability [?; Pearce 1984]. We

repeat the definition here, using the notation of Osborne and Rubinstein [1994].

**Definition 2.1:** A strategy  $\sigma$  for player  $i$  in game  $\Gamma$  is *rationalizable* if, for each player  $j$ , there exists a sequence  $X_j^0, X_j^1, X_j^2, \dots$  of sets of strategies for player  $j$  such that  $X_j^0 = \Sigma_j$  and, for each strategy  $\sigma' \in X_j^k$ ,  $k \geq 1$ , a probability measure  $\mu_{\sigma', k}$  whose support is a subset of  $\bar{X}_j^{k-1}$  such that  $\sigma \in \bigcap_{j=0}^{\infty} X_j$  and, for each player  $j$ , each strategy  $\sigma' \in X_j^k$  is a best response to the beliefs  $\mu_{\sigma', k}$ . ■

Intuitively,  $X_j^1$  consists of strategies that are best responses to some belief of player  $j$ , and  $X_j^{h+1}$  consists of strategies in  $X_j^h$  that are best responses to some belief of player  $j$  with support  $X_j^h$ ; that is, beliefs that assume that everyone else is best responding to some beliefs assuming that everyone else is responding to some beliefs assuming  $\dots$  ( $h$  times).

We now state the epistemic characterization of rationalizability due to Tan and Werlang [1988] in our language; it just says that a strategy is rationalizable iff it can be played in a state where rationality is common knowledge.

Let  $RAT$  be an abbreviation for  $RAT_1 \wedge \dots \wedge RAT_n$ ; let  $E\varphi$  be an abbreviation of  $B_1\varphi \wedge \dots \wedge B_n\varphi$ ; and define  $E^k\varphi$  for all  $k$  inductively by taking  $E^0\varphi$  to be  $\varphi$  and  $E^{k+1}\varphi$  to be  $E(E^k\varphi)$ . Common knowledge of  $\varphi$  holds iff  $E^k\varphi$  holds for all  $k \geq 0$ .

**Theorem 2.2:** *The following are equivalent:*

- (a)  $\sigma$  is a rationalizable strategy for  $i$  in game  $\Gamma$ ;
- (b) there exists a measurable structure  $M$  that is appropriate for  $\Gamma$  and a state  $\omega$  such that  $\mathbf{s}_i(\omega) = \sigma$  and  $(M, \omega) \models B_i E^k RAT$  for all  $k \geq 0$ ;
- (c) there exists a structure  $M$  that is appropriate for  $\Gamma$  and a state  $\omega$  such that  $\mathbf{s}_i(\omega) = \sigma$  and  $(M, \omega) \models B_i E^k RAT$  for all  $k \geq 0$ .

Since the proof is similar in spirit to that of Tan and Werlang [1988], we omit it here.

We now consider iterated deletion of strongly dominated (resp., weakly dominated) strategies.

**Definition 2.3:** Strategy  $\sigma$  for player  $i$  is *strongly dominated* by  $\sigma'$  with respect to  $\Sigma'_{-i} \subseteq \Sigma_{-i}$  if  $u_i(\sigma, \tau_{-i}) > u_i(\sigma', \tau_{-i})$  for all  $\tau_{-i} \in \Sigma'_{-i}$ . Strategy  $\sigma$  for player  $i$  is *weakly dominated* by  $\sigma'$  with respect to  $\Sigma'_{-i} \subseteq \Sigma_{-i}$

if  $u_i(\sigma, \tau_{-i}) \geq u_i(\sigma, \tau'_{-i})$  for all  $\tau_{-i} \in \Sigma'_{-i}$  and  $u_i(\sigma, \tau'_{-i}) > u_i(\sigma, \tau''_{-i})$  for some  $\tau''_{-i} \in \Sigma'_{-i}$ .

Strategy  $\sigma$  for player  $i$  survives  $k$  rounds of iterated deletion of strongly dominated (resp., weakly dominated) strategies if, for each player  $j$ , there exists a sequence  $X_j^0, X_j^1, X_j^2, \dots, X_j^k$  of sets of strategies for player  $j$  such that  $X_j^0 = \Sigma_j$  and, if  $h < k$ , then  $X_j^{h+1}$  consists of the strategies in  $X_j^h$  not strongly (resp., weakly) dominated by any strategy with respect to  $X_{-j}^h$ , and  $\sigma \in X_i^k$ . Strategy  $\sigma$  survives iterated deletion of strongly dominated (resp., weakly dominated) strategies if it survives  $k$  rounds of iterated deletion for all  $k$ . ■

The following well-known result connects strong and weak dominance to best responses.

**Proposition 2.4:** [Pearce 1984]

- A strategy  $\sigma$  for player  $i$  is not strongly dominated by any strategy with respect to  $\Sigma'_{-i}$  iff there is a belief  $\mu_\sigma$  of player  $i$  whose support is a subset of  $\Sigma'_{-i}$  such that  $\sigma$  is a best response with respect to  $\mu_\sigma$ .
- A strategy  $\sigma$  for player  $i$  is not weakly dominated by any strategy with respect to  $\Sigma'_{-i}$  iff there is a belief  $\mu_\sigma$  of player  $i$  whose support is all of  $\Sigma'_{-i}$  such that  $\sigma$  is a best response with respect to  $\mu_\sigma$ .

It immediately follows from Proposition 2.4 (and is well known) that a strategy is rationalizable iff it survives iterated deletion of strongly dominated strategies. Thus, the characterization of rationalizability in Theorem 2.2 is also a characterization of strategies that survive iterated deletion of strongly dominated strategies. We now give a slightly different characterization that allows us to relate iterated deletion of strongly and weakly dominated strategies. For each player  $i$ , define the formulas  $C_i^k$  inductively by taking  $C_i^0$  to be *true* and  $C_i^{k+1}$  to be an abbreviation of

$$RAT_i \wedge B_i(\bigwedge_j C_j^k).$$

That is,  $C_i^{k+1}$  holds (i.e., player  $i$  is  $k+1$ -level rational) iff player  $i$  is playing a best response to his beliefs, and he knows that all players are  $k$ -level rational.

Because players know their own strategies and beliefs, it easily follows that  $RAT_i \Rightarrow B_i RAT_i$  is valid (i.e., true in all states in all appropriate models). Two formulas  $\varphi$  and  $\psi$  are *logically equivalent* if  $\varphi \Leftrightarrow \psi$  is valid.

**Lemma 2.5:** *The following formulas are logically equivalent for all  $k \geq 0$ .*

- $C_i^{k+1}$ ;
- $B_i C_i^{k+1}$ ;
- $RAT_i \wedge B_i(\bigwedge_{j \neq i} C_j^k)$ ;
- $RAT_i \wedge B_i E^{k-1} RAT$  (taking  $E^{-1}\varphi$  to be true).

Moreover,  $C_i^{k+1} \Rightarrow C_i^k$  is valid for all  $k \geq 0$ .

**Proof:** A straightforward induction on  $k$ . ■

The following alternative characterization of iterated deletion of strongly dominated strategies follow in the same way as Theorem 2.2.

**Theorem 2.6:** *The following are equivalent:*

- the strategy  $\sigma$  for player  $i$  survives  $k$  rounds of iterated deletion of strongly dominated strategies in game  $\Gamma$ ;
- there is a measurable structure  $M^k$  appropriate for  $\Gamma$  and a state  $\omega^k$  in  $M^k$  such that  $\mathbf{s}_i(\omega^k) = \sigma$  and  $(M^k, \omega^k) \models C_i^k$ ;
- there is a structure  $M^k$  appropriate for  $\Gamma$  and a state  $\omega^k$  in  $M^k$  such that  $\mathbf{s}_i(\omega^k) = \sigma$  and  $(M^k, \omega^k) \models C_i^k$ .

The following corollary now follows from Theorem 2.6, Lemma 2.5, and the fact that the deletion procedure converges after a finite number of steps.

**Corollary 2.7:** *The following are equivalent:*

- The strategy  $\sigma$  for player  $i$  survives iterated deletion of strongly dominated strategies in game  $\Gamma$ ;
- there exists a measurable structure  $M$  that is appropriate for  $\Gamma$  and a state  $\omega$  such that  $\mathbf{s}_i(\omega) = \sigma$  and  $(M, \omega) \models C_i^k$  for all  $k \geq 0$ ;
- there exists a structure  $M$  that is appropriate for  $\Gamma$  and a state  $\omega$  such that  $\mathbf{s}_i(\omega) = \sigma$  and  $(M, \omega) \models C_i^k$  for all  $k \geq 0$ ;

We next turn to characterizing iterated deletion of weakly dominated strategies.

### 3 Characterizing Iterated Admissibility

In the standard treatment (which is essentially the one considered in Section 2), player  $i$  is taken to be  $(k+1)$ -level rational iff player  $i$  is rational (i.e., playing a best response to his beliefs), and knows that all other player are  $k$ -level rational.<sup>2</sup> But what else do players know?

We want to consider a situation where, intuitively, *all* an agent knows about the other agents is that they satisfy the appropriate rationality assumptions. More precisely, we modify the formula  $C_i^{k+1}$  to require that not only does player  $i$  know that the players are  $k$ -level rational, but this is the *only* thing that he knows about the other players. That is, we say that agent  $i$  is  $(k+1)$ -level rational if player  $i$  is rational, he knows that the players are  $k$ -level rational, and this is all player  $i$  knows about the other players. We here use the phrase “all agent  $i$  knows” in essentially the same sense that it is used by Levesque [1990] and Halpern and Lakemeyer [2001], but formalize it a bit differently. Roughly speaking, we interpret “all agent  $i$  knows is  $\varphi$ ” as meaning that agent  $i$  believes  $\varphi$ , and considers possible every *formula* about the other players that is consistent with  $\varphi$ . Thus, what “all I know” means is very sensitive to the choice of the language.

To formalize this, consider the modal operator  $\diamond$  defined as follows:

- $(M, \omega) \models \diamond\varphi$  iff there is some structure  $M'$  appropriate for  $\Gamma$  and state  $\omega'$  such that  $(M', \omega') \models \varphi$ .

Intuitively,  $\diamond\varphi$  is true if there is some state and structure where  $\varphi$  is true; that is, if  $\varphi$  is satisfiable. Note that if  $\diamond\varphi$  is true at some state, then it is true at all states in all structures. Define  $O_i^{\mathcal{L}}\varphi$  (read “all agent  $i$  knows with respect to the language  $\mathcal{L}$ ”) to be an abbreviation for

$$B_i\varphi \wedge (\wedge_{\psi \in \mathcal{L}} (\diamond(\varphi \wedge \psi) \Rightarrow \langle B_i \rangle \psi)).$$

Since  $O_i^{\emptyset}$  (where  $\emptyset$  denotes the empty language) is equivalent to  $B_i$  (under the standard identification of the empty conjunction with the formula *true*), it follows that  $C_i^{k+1}$  is just  $RAT_i \wedge O_i^{\emptyset}(\wedge_j C_j^k)$ . We next consider a slightly richer language, whose formulas

<sup>2</sup>We should perhaps say “believes” here rather than “knows”, since a player can be mistaken. We are deliberately blurring the subtle distinctions between “knowledge” and “belief” here.

can talk about strategies (but not beliefs) of the players. (In Section 4, we consider languages in which we can talk about both the strategies and beliefs of the players.) Define the primitive proposition  $play_i(\sigma)$  as follows:

- $(M, \omega) \models play_i(\sigma)$  iff  $\omega \in \llbracket \sigma \rrbracket_M$ .

Let  $play(\vec{\sigma})$  be an abbreviation for  $\wedge_{j=1}^n play_j(\sigma_j)$ , and let  $play_{-i}(\sigma_{-i})$  be an abbreviation for  $\wedge_{j \neq i} play_j(\sigma_j)$ . Intuitively,  $(M, \omega) \models play(\vec{\sigma})$  iff  $s(\omega) = \sigma$ , and  $(M, \omega) \models play_{-i}(\sigma_{-i})$  if, at  $\omega$ , the players other than  $i$  are playing strategy profile  $\sigma_{-i}$ .

Let  $\mathcal{L}^0(\Gamma)$  be the language whose only formulas are (Boolean combinations of) formulas of the form  $play_i(\sigma)$ ,  $i = 1, \dots, n$ ,  $\sigma \in \Sigma_i$ . Let  $\mathcal{L}_i^0(\Gamma)$  consist of just the formulas of the form  $play_i(\sigma)$ , and let  $\mathcal{L}_{-i}^0(\Gamma) = \cup_{j \neq i} \mathcal{L}_j^0(\Gamma)$ . Again, we omit the parenthetical  $\Gamma$  when it is clear from context or irrelevant.

We would now like to define  $D_i^{k+1}$  by replacing  $B_i$  by  $O_i^{\mathcal{L}_{-i}^0}$  in the definition of  $C_i^{k+1}$ . However, the resulting formulas are in general inconsistent! For example,  $D_i^3$  would then imply  $D_i^2$ , which would require that  $i$  consider possible (i.e., assign positive probability to) all strategies for the other players consistent with 1-rationality, while  $D_i^3$  would also require that player  $i$  believe that the other players are 2-rational (and thus player  $i$  should ascribe probability 0 to all strategies for the other players that are 1-rational but not 2-rational). The first step in getting an appropriate analogue to  $C_i^k$  is to remove the conjunct  $D_i^k$  from the scope of  $O_i^{\mathcal{L}_{-i}^0}$ . As Lemma 2.5(c) shows, this change would have no impact on the definition of  $C_i^{k+1}$ . With this change, it is no longer the case that  $D_i^{k+1}$  implies  $D_i^k$  (which is a good thing, since the two formulas are still inconsistent). However, since we want  $D_i^{k+1}$  to be true at a state if the strategy used by player  $i$  at that state survives  $k$  rounds of iterated deletion of weakly dominated strategies, we need to add a conjunct to  $D_i^{k+1}$  that guarantees that the strategy used is  $k'$ -rational for  $k' < k+1$ . Thus, for each player  $i$ , define the formulas  $D_i^k$  inductively by taking  $D_i^0$  to be *true* and  $D_i^{k+1}$  to be an abbreviation of

$$RAT_i \wedge D_i^k \wedge O_i^{\mathcal{L}_{-i}^0}(\wedge_{j \neq i} D_j^k).$$

where  $D_i^k$  (read “player  $i$  plays a strategy consistent with  $k$ -level rationality”) is an abbreviation of  $(play_i(\sigma) \Rightarrow \diamond(play_i(\sigma) \wedge D_i^k))$ . That is,  $D_i^{k+1}$  holds (i.e., player  $i$  is  $(k+1)$ -level rational) iff player  $i$  is

rational, plays a strategy that is consistent with  $k$ -level rationality, knows that other players are  $k$ -level rational, and that is all player  $i$  knows about the *strategies* of the other players.<sup>3</sup>

By expanding the modal operator  $O$ , it easily follows that  $D_i^{k+1}$  implies  $RAT_i \wedge B_i(\wedge_{j \neq i} D_j^k)$ ; an easy induction on  $k$  then shows that  $D_i^{k+1}$  implies  $C_i^{k+1}$ . But  $D_i^{k+1}$  requires more; it requires that player  $i$  assigns positive probability to each strategy profile for the other players that is compatible with  $D_{-i}^k$  (i.e., with level- $k$  rationality). As we now show, the formula  $D_i^k$  characterizes strategies that survive iterated deletion of weakly dominated strategies.

**Theorem 3.1:** *The following are equivalent:*

- (a) *the strategy  $\sigma$  for player  $i$  survives  $k$  rounds of iterated deletion of weakly dominated strategies;*
- (b) *there is a measurable structure  $M^k$  appropriate for  $\Gamma$  and a state  $\omega^k$  in  $M^k$  such that  $\mathbf{s}_i(\omega^k) = \sigma$  and  $(M^k, \omega^k) \models D_i^k$ ;*
- (c) *there is a structure  $M^k$  appropriate for  $\Gamma$  and a state  $\omega^k$  in  $M^k$  such that  $\mathbf{s}_i(\omega^k) = \sigma$  and  $(M^k, \omega^k) \models D_i^k$ .*

*In addition, there is a finite structure  $\overline{M}^k = (\Omega^k, \mathbf{s}, \mathcal{F}, \mathcal{PR}_1, \dots, \mathcal{PR}_n)$  such that  $\Omega^k = \{(k', i, \vec{\sigma}) : k' \leq k, 1 \leq i \leq n, \vec{\sigma} \in X_1^{k'} \times \dots \times X_n^{k'}\}$ ,  $\mathbf{s}(k', i, \vec{\sigma}) = \vec{\sigma}$ ,  $\mathcal{F} = 2^{\Omega^k}$ , where  $X_j^{k'}$  consists of all strategies for player  $j$  that survive  $k'$  rounds of iterated deletion of weakly dominated strategies and, for all states  $(k', i, \vec{\sigma}) \in \Omega^k$ ,  $(\overline{M}^k, (k', i, \vec{\sigma})) \models \wedge_{j \neq i} D_j^{k'}$ .*

**Proof:** We proceed by induction on  $k$ , proving both the equivalence of (a), (b), and (c) and the existence of a structure  $\overline{M}^k$  with the required properties.

The result clearly holds if  $k = 0$ . Suppose that the result holds for  $k$ ; we show that it holds for  $k + 1$ . We first show that (c) implies (a). Suppose that  $(M^{k+1}, \omega^{k+1}) \models D_i^{k+1}$  and  $\mathbf{s}_i(\omega^{k+1}) = \sigma_i$ . It follows that  $\sigma_i$  is a best response to the belief  $\mu_{\sigma_i}$  on the strategies of other players induced by  $\mathcal{PR}_i^{k+1}(\omega)$ . Since  $(M^{k+1}, \omega^{k+1}) \models B_i(\wedge_{j \neq i} D_j^k)$ , it follows from the induction hypothesis that the support of  $\mu_{\sigma_i}$  is contained in  $X_{-i}^k$ .

<sup>3</sup>  $D_i^k$  is essentially the formula  $RAT_i^k$  from the introduction. However, since we consider a number of variants of  $RAT_i^k$ , we find it useful to distinguish them.

Since  $(M^{k+1}, \omega^{k+1}) \models \wedge_{\sigma_{-i} \in \Sigma_{-i}} (\diamond(\text{play}_{-i}(\sigma_{-i}) \wedge (\wedge_{j \neq i} D_j^k))) \Rightarrow \langle B_j \rangle(\text{play}_{-i}(\sigma_{-i}))$ , it follows from the induction hypothesis that the support of  $\mu_{\sigma_i}$  is all of  $X_{-i}^k$ . Since  $(M^{k+1}, \omega^{k+1}) \models D_i^k$ , it follows from the induction hypothesis that  $\sigma_i \in X_i^k$ . Thus, since  $(M^{k+1}, \omega^{k+1}) \models RAT_i$ , it follows by Proposition 2.4 that  $\sigma_i \in X_i^{k+1}$ .

We next construct the structure  $\overline{M}^{k+1} = (\Omega^{k+1}, \mathbf{s}, \mathcal{F}, \mathcal{PR}_1, \dots, \mathcal{PR}_n)$ . As required, we define  $\Omega^{k+1} = \{(k', i, \vec{\sigma}) : k' \leq k + 1, 1 \leq i \leq n, \vec{\sigma} \in X_1^{k'} \times \dots \times X_n^{k'}\}$ ,  $\mathbf{s}(k', i, \vec{\sigma}) = \vec{\sigma}$ ,  $\mathcal{F} = 2^{\Omega^{k+1}}$ . For a state  $\omega$  of the form  $(k', i, \vec{\sigma})$ , since  $\sigma_j \in X_j^{k'}$ , by Proposition 2.4, there exists a distribution  $\mu_{k', \sigma_j}$  whose support is all of  $X_{-j}^{k-1}$  such that  $\sigma_j$  is a best response to  $\mu_{\sigma_j}$ . Extend  $\mu_{k', \sigma_j}$  to a distribution  $\mu_{k', i, \sigma_j}$  on  $\Omega^{k+1}$  as follows:

- for  $i \neq j$ , let  $\mu_{k', i, \sigma_j}^{k', i, \sigma_j}(k'', i', \vec{\tau}) = \mu_{k', \sigma_j}(\vec{\tau}_{-j})$  if  $i' = j, k'' = k' - 1$ , and  $\tau_j = \sigma_j$ , and 0 otherwise;
- $\mu_{k', i, \sigma_j}^{k', i, \sigma_j}(k'', i', \vec{\tau}) = \mu_{k', \sigma_j}(\vec{\tau}_{-j})$  if  $i' = j, k'' = k'$ , and  $\tau_j = \sigma_j$ , and 0 otherwise.

Let  $\mathcal{PR}_j(k', i, \vec{\sigma}) = \sigma_{k', i, \sigma_j}^{k', i, \sigma_j}$ . We leave it to the reader to check that this structure is appropriate. An easy induction on  $k'$  now shows that  $(\overline{M}^{k+1}, (k', i, \vec{\sigma})) \models \wedge_{j \neq i} D_j^{k'}$  for  $i = 1, \dots, n$ .

To see that (a) implies (b), suppose that  $\sigma_j \in X_j^{k+1}$ . Choose a state  $\omega$  in  $\overline{M}^{k+1}$  of the form  $(k + 1, i, \vec{\sigma})$ , where  $i \neq j$ . As we just showed,  $(\overline{M}^{k+1}, (k', i, \vec{\sigma})) \models D_j^{k'}$ , and  $\mathbf{s}_j(k', i, \vec{\sigma}) = \sigma_j$ . Moreover,  $\overline{M}^{k+1}$  is measurable (since  $\mathcal{F}$  consists of all subsets of  $\Omega^{k+1}$ ).

Clearly (b) implies (c). ■

Note that there is no analogue of Corollary 2.7 here. This is because there is no state where  $D_i^k$  holds for all  $k \geq 0$ ; it cannot be the case that  $i$  places positive probability on all strategies (as required by  $D_1^k$ ) and that  $i$  places positive probability only on strategies that survive one round of iterated deletion (as required by  $D_2^k$ ), unless all strategies survive one round on iterated deletion. We can say something slightly weaker though. There is some  $k^*$  such that the process of iterated deletion converges after  $k^*$  steps; that is,  $X_j^{k^*} = X_j^{k^*+1}$  for all  $j$  (and hence  $X_j^{k^*} = X_j^{k'}$  for all  $k' \geq k^*$ ). That means that there is a state where  $D_i^{k'}$  holds for all  $k' > k^*$ . Thus, we can show that a strategy  $\sigma$  for player  $i$  survives iterated deletion of weakly dominated strategies iff there exists a  $k^*$  and

a state  $\omega$  such that  $s_i(\omega) = \sigma$  and  $(M, \omega) \models D_i^{k'}$  for all  $k' > k^*$ . Since  $C_i^{k+1}$  implies  $C_i^k$ , an analogous result holds for iterated deletion of strongly dominated strategies, with  $D_i^{k'}$  replaced by  $C_i^{k'}$ .

It is also worth stressing that, unlike the BFK construction, in a state where  $D^k$  holds, an agent does *not* consider all strategies possible, but only the ones consistent with the appropriate level of rationality. We could require the agent to consider all strategies possible by using LPSs or nonstandard probability. The only change that this would make to our characterization is that, if we are using nonstandard probability, we would interpret  $B_i\varphi$  to mean that  $\varphi$  holds with probability infinitesimally close to 1, while  $\langle B_i \rangle \varphi$  would mean that  $\varphi$  holds with probability whose standard part is positive (i.e., non-infinitesimal probability). We do not pursue this point further.

## 4 Richer Languages

The formula  $D_i^{k+1}$  was defined with respect to the language  $\mathcal{L}_{-i}^O$ , and thus required that player  $i$  assign positive probability to all and only strategies consistent with  $D_{-i}^k$ . But why focus just on strategies? We now consider richer languages that can talk about players' beliefs; this requires players to ascribe positive probability to all beliefs that the other agents could have as well as all the strategies they could be using (which are consistent with appropriate levels of rationality).

First, let  $\mathcal{L}^2(\Gamma)$  be the extension of  $\mathcal{L}^1$  that includes a primitive proposition  $play_i(\sigma)$  for each player  $i$  and strategy  $\sigma \in \Sigma_i$ .

To relate our results to those of BFK, even the language  $\mathcal{L}^2$  is too weak, since it does not allow an agent to express probabilistic beliefs. Let  $\mathcal{L}^3(\Gamma)$  be the language that extends  $\mathcal{L}^2(\Gamma)$  by allowing formulas of the form  $pr_i(\varphi) \geq \alpha$  and  $pr_i(\varphi) > \alpha$ , where  $\alpha$  is a rational number in  $[0, 1]$ ;  $pr_i(\varphi) \geq \alpha$  can be read as “the probability of  $\varphi$  according to  $i$  is at least  $\alpha$ ”, and similarly for  $pr_i(\varphi) > \alpha$ . We allow nesting here, so that we can have a formula of the form  $pr_j(play_i(\sigma) \wedge pr_k(play_i(\sigma')) > 1/3) \geq 1/4$ . As we would expect,

- $(M, \omega) \models pr_i(\varphi)$  iff  $\mathcal{PR}_i(\omega)(\llbracket \varphi \rrbracket_M) \geq \alpha$ .

The restriction to  $\alpha$  being rational allows the language to be countable. However, as we now show, it is not too serious a restriction.

Let  $\mathcal{L}^4(\Gamma)$  be the language that extends  $\mathcal{L}^2(\Gamma)$  by closing off under countable conjunctions, so that if  $\varphi_1, \varphi_2, \dots$  are formulas, then so is  $\bigwedge_{m=1}^{\infty} \varphi_m$ , and formulas of the form  $pr_i(\varphi) > \alpha$ , where  $\alpha$  is a real number in  $[0, 1]$ . (We can express  $pr_i(\varphi) \geq \alpha$  as the countable conjunction  $\bigwedge_{\beta < \alpha, \beta \in Q \cap [0, 1]} pr_i(\varphi) > \beta$ , where  $Q$  is the set of rational numbers, so there is no need to include formulas of the form  $pr_i(\varphi) \geq \alpha$  explicitly in  $\mathcal{L}^4(\Gamma)$ .) We omit the parenthetical  $\Gamma$  in  $\mathcal{L}^3(\Gamma)$  and  $\mathcal{L}^4(\Gamma)$  when the game  $\Gamma$  is clear from context. A subset  $\Phi$  of  $\mathcal{L}^3$  is  $\mathcal{L}^3$ -realizable if there exists an appropriate structure  $M$  for  $\Gamma$  and state  $\omega$  in  $M$  such that, for all formulas  $\varphi \in \mathcal{L}^3$ ,  $(M, \omega) \models \varphi$  iff  $\varphi \in \Phi$ .<sup>4</sup> We can similarly define what it means for a subset of  $\mathcal{L}^4$  to be  $\mathcal{L}^4$ -realizable.

The following lemma is proved in the full paper.

**Lemma 4.1:** *Every  $\mathcal{L}^3$ -realizable set can be uniquely extended to an  $\mathcal{L}^4$ -realizable set.*

With this background, let  $\mathcal{L}_i^3$  consist of all formulas in  $\mathcal{L}^3$  of the form  $pr_i(\varphi) \geq \alpha$  and  $pr_i(\varphi) > \alpha$  ( $\varphi$  can mention  $pr_j$ ,  $j \neq i$ ; it is only the outermost modal operator that must be  $i$ ). Intuitively,  $\mathcal{L}_i^3$  consists of the formulas describing  $i$ 's beliefs. Let  $\mathcal{L}_{i+}^3$  consist of  $\mathcal{L}_i^3$  together with formulas of the form  $true$ ,  $RAT_i$ , and  $play_i(\sigma)$ , for  $\sigma \in \Sigma_i$ . Let  $\mathcal{L}_{(-i)+}^3$  be an abbreviation for  $\bigcup_{j \neq i} \mathcal{L}_{j+}^3$ . We can similarly define  $\mathcal{L}_i^4$  and  $\mathcal{L}_{i+}^4$ .

Note that if  $\varphi \in \mathcal{L}_{(-i)+}^3$ , then  $O_i^{L^3(-i)+} \varphi$  is an abbreviation for the formula

$$B_i\varphi \wedge (\bigwedge_{\psi \in \mathcal{L}_{(-i)+}^3} \diamond(\varphi \wedge \psi)) \Rightarrow \langle B_j \rangle \psi.$$

Thus,  $O_i^{L^3(-i)+} \varphi$  holds if agent  $i$  believes  $\varphi$  but does not know anything beyond that; he ascribes positive probability to all formulas in  $\mathcal{L}_{(-i)+}^3$  consistent with  $\varphi$ . This is very much in the spirit of the Halpern-Lakemeyer [2001] definition in the context of epistemic logic. Of course, we could go further and define a notion of “all  $i$  knows” for the language  $\mathcal{L}^4$ . Doing this would give a definition that is even closer to that of Halpern and Lakemeyer. Unfortunately, we cannot require than agent  $i$  ascribe positive probability to all the formulas in  $\mathcal{L}_{(-i)+}^4$  consistent with  $\varphi$ ; in general, there will be an uncountable number of distinct and mutually exclusive formulas consistent with  $\varphi$ , so they cannot all be assigned positive probability. This problem

<sup>4</sup>For readers familiar with standard completeness proofs in modal logic, if we had axiomatized the logic we are implicitly using here, the  $\mathcal{L}^3$ -realizable sets would just be the maximal consistent sets in the logic.

does not arise with  $\mathcal{L}^3$ , since it is a countable language. Halpern and Lakemeyer could allow an agent to consider an uncountable set of worlds possible, since they were not dealing with probabilistic systems. In the sequel, we thus focus on the language  $\mathcal{L}^3$  and let  $O_i$  denote  $O_i^{\mathcal{L}^3(-i)+}$ .<sup>5</sup>

Define the formulas  $F_i^k$  inductively by taking  $F_i^0$  to be the formula *true*, and  $F_i^{k+1}$  to an abbreviation of

$$RAT_i \wedge F_i^k \wedge O_i(\bigwedge_{j \neq i} F_j^k),$$

where  $F_i^k$  is an abbreviation of  $(play_i(\sigma) \Rightarrow \diamond(play_i(\sigma) \wedge F_i^k))$ . Thus,  $F_i^{k+1}$  says that  $i$  is rational, plays a  $k$ -level rational strategy, knows that all the other players satisfy level- $k$  rationality (i.e.,  $F_j^k$ ), and that is all that  $i$  knows. It is easy to see that  $F_j^{k+1}$  implies  $D_j^{k+1}$ . The difference is that instead of requiring just that  $j$  assign positive probability to all strategy profiles compatible with  $F_{-j}^k$ , it requires that  $j$  assign positive probability to all formulas in  $\mathcal{L}_{(-i)+}^3$  compatible with  $F_{-j}^k$ .

The next result shows that  $F_i^k$  characterizes iterated admissibility, just as  $D_i^k$  does.

**Theorem 4.2:** *The following are equivalent:*

- (a) *the strategy  $\sigma$  for player  $i$  survives  $k$  rounds of iterated deletion of weakly dominated strategies;*
- (b) *there is a measurable structure  $M^k$  appropriate for  $\Gamma$  and a state  $\omega^k$  in  $M^k$  such that  $s_i(\omega^k) = \sigma$  and  $(M^k, \omega^k) \models F_i^k$ ;*
- (c) *there is a structure  $M^k$  appropriate for  $\Gamma$  and a state  $\omega^k$  in  $M^k$  such that  $s_i(\omega^k) = \sigma$  and  $(M^k, \omega^k) \models F_i^k$ ;*

The proof of Theorem 4.2 is similar in spirit to the proof of Theorem 3.1, but is more complicated. Due to lack of space, we defer it to the full version.

## 5 Complete and Canonical Structures

### 5.1 Canonical Structures

The intuition behind “all  $i$  knows is  $\varphi$ ” goes back to Levesque [1990]. The idea is that all  $i$  knows is  $\varphi$  if

<sup>5</sup>Note that the modal operator  $\diamond$  is not in the language  $\mathcal{L}^3$  or  $\mathcal{L}^4$ . None of our results would be affected if we had considered a language that also included  $\diamond$ ; for ease of exposition, we have decided not to include  $\diamond$  here.

(1)  $i$  knows  $\varphi$  (so that  $\varphi$  is true in all the worlds that  $i$  considers possible) and (2)  $i$  considers possible all worlds consistent with  $\varphi$  (so that  $\varphi$  is false in all worlds that  $i$  does not consider possible). In the single-agent case (which is what Levesque considered) it is relatively easy to make precise the set of worlds that  $i$  does not consider possible, since a world can be identified with a truth assignment. This is much more complicated in the multi-agent setting considered by Halpern and Lakemeyer [2001]. They made it precise by working in the canonical structure. The advantage of the canonical structure is that, in a sense it has all possible worlds, so it is clear what worlds an agent does not consider possible. Although our definition of “all  $i$  knows” is more language-dependent, our intuitions for the notion are still grounded in the canonical structure. Thus, in this section, we consider our definitions in the context of canonical structures. The reason we say “structures” here rather than “structure” is that the notion of canonical structure is also language dependent.

Define the *canonical structure*  $M^c = (\Omega^c, s^c, \mathcal{F}^c, \mathcal{PR}_1^c, \dots, \mathcal{PR}_n^c)$  for  $\mathcal{L}^4$  as follows:

- $\Omega^c = \{\omega_\Phi : \Phi \text{ is a realizable subset of } \mathcal{L}^4(\Gamma)\};$
- $s^c(\omega_\Phi) = \vec{\sigma}$  iff  $play(\sigma) \in \Phi$ ;
- $\mathcal{F}^c = \{F_\varphi : \varphi \in \mathcal{L}^4\}$ , where  $F_\varphi = \{\omega_\Phi : \varphi \in \Phi\}$ ;
- $\text{Pr}_i^c(\omega_\Phi)(F_\varphi) = \inf\{\alpha : \text{pr}_i(\varphi) > \alpha \in \Phi\}$ .

**Lemma 5.1:**  *$M^c$  is an appropriate measurable structure for  $\Gamma$ .*

The following result is the analogue of the standard “truth lemma” in completeness proofs in modal logic.

**Proposition 5.2:** *For  $\psi \in \mathcal{L}^4$ ,  $(M^c, \omega_\Phi) \models \psi$  iff  $\psi \in \Phi$ .*

We have constructed a canonical structure for  $\mathcal{L}^4$ . It follows easily from Lemma 4.1 that the canonical structure for  $\mathcal{L}^3$  (where the states are realizable  $\mathcal{L}^3$  sets) is isomorphic to  $M^c$ . (In this case, the set  $\mathcal{F}^c$  of measurable sets would be the smallest  $\sigma$ -algebra containing  $\llbracket \varphi \rrbracket_M$  for  $\varphi \in \mathcal{L}^3$ .) Thus, the choice of  $\mathcal{L}^3$  vs.  $\mathcal{L}^4$  does not play an important role when constructing a canonical structure.

A strategy  $\sigma_i$  for player  $i$  survives iterated deletion of weakly dominated strategies iff the formula

$$\text{undominated}^k(\sigma_i) = play_i(\sigma_i) \wedge \exists k(\bigwedge_{k=k^*}^\infty F_i^k)$$



is satisfied at some state in the canonical structure. But there are other structures in which  $undominated(\sigma_i)$  is satisfied. One way to get such a structure is by essentially “duplicating” states in the canonical structure. The canonical structure can be *embedded* in a structure  $M$  if, for all  $\mathcal{L}^3$ -realizable sets  $\Phi$ , there is a state  $\omega_\Phi$  in  $M$  such that  $(M, \omega_\Phi) \models \varphi$  iff  $\varphi \in \Phi$ . Clearly  $undominated(\sigma_i)$  is satisfied in any structure in which the canonical structure can be embedded.

A structure in which the canonical structure can be embedded is in a sense larger than the canonical structure. But  $undominated(\sigma_i)$  can be satisfied in structures smaller than the canonical structure. There are two reasons for this. The first is that to satisfy  $undominated(\sigma_i)$ , there is no need to consider a structure with states where all the players are irrational. It suffices to restrict to states where at least one player is using a strategy that survives at least one round of iterated deletion. This is because players know their strategy; thus, in a state where a strategy  $\sigma_i$  for player  $i$  is admissible, player  $i$  must ascribe positive probability to all other strategies; however, in those states, player  $i$  still plays  $\sigma_i$ .

A perhaps more interesting reason that we do not need the canonical structure is our use of the language  $\mathcal{L}_3$ . The formulas  $F_i^k$  guarantee that player  $i$  ascribes positive probability to all formulas  $\varphi$  consistent with the appropriate level of rationality. Since a finite conjunction of formulas in  $\mathcal{L}^3$  is also a formula in  $\mathcal{L}^3$ , player  $i$  will ascribe positive probability to all finite conjunctions of formulas consistent with rationality. But a state is characterized by a *countable* conjunction of formulas. Since  $\mathcal{L}^3$  is not closed under countable conjunctions, a structure that satisfies  $undominated(\sigma_i)$  may not have states corresponding to all  $\mathcal{L}^3$ -realizable sets of formulas. If we had used  $\mathcal{L}^4$  instead of  $\mathcal{L}^3$  in the definition of  $F_i^k$  (ignoring the issues raised earlier with using  $\mathcal{L}^4$ ), then there would be a state corresponding to every  $\mathcal{L}^4$ -realizable (equivalently,  $\mathcal{L}^3$ -realizable) set of formulas. Alternatively, if we consider appropriate structures that are compact in a topology where all sets definable by formulas (i.e., sets of the form  $\llbracket \varphi \rrbracket_M$ , for  $\varphi \in \mathcal{L}^3$ ) are closed (in which case they are also open, since  $\llbracket \neg \varphi \rrbracket_M$  is the complement of  $\llbracket \varphi \rrbracket_M$ ), then all states where at least one player is using a strategy that survives at least one round of iterated deletion will be in the structure.

Although, as this discussion makes clear, the formula  $F_i^k$  that characterizes iterated admissibility can be sat-

isfied in structures quite different from the canonical structure, the canonical structure does seem to be the most appropriate setting for reasoning about statements involving “all agent  $i$  knows”. Moreover, as we now show, canonical structures allow us to relate our approach to that of BFK.

## 5.2 Complete Structures

BFK worked with complete structures. We now want to show that  $M^c$  is complete, in the sense of BFK. To make this precise, we need to recall some notions from BFK (with some minor changes to be more consistent with our notation).

BFK considered what they called *interactive probability structures*. These can be viewed as a special case of probability structures. A *BFK-like structure* (for a game  $\Gamma$ ) is a probability structure  $M = (\Omega, \mathfrak{s}, \mathcal{F}, \mathcal{PR}_1, \dots, \mathcal{PR}_n)$  such that there exist spaces  $T_1, \dots, T_n$  (where  $T_i$  can be thought of as the *type space* for player  $i$ ) such that

- $\Omega$  is isomorphic to  $\vec{\Sigma} \times \vec{T}$  via some isomorphism  $h$ ;
- if  $h(\omega) = \vec{\sigma} \times \vec{t}$ , then
  - $\mathfrak{s}(\omega) = \vec{\sigma}$ ;
  - taking  $T_i(\omega) = t_i$  (i.e., the type of player  $i$  in  $h(\omega)$  is  $t_i$ ); the support of  $\mathcal{PR}_i(\omega)$  is contained in  $\{\omega' : \mathfrak{s}_i(\omega') = \sigma', T_i(\omega') = t_i\}$ , so that  $\mathcal{PR}_i(\omega)$  induces a probability on  $\Sigma_{-i} \times T_{-i}$ ;
  - $\mathcal{PR}_i(\omega)$  depends only on  $T_i(\omega)$ , in the sense that if  $T_i(\omega) = T_i(\omega')$ , then  $\mathcal{PR}_i(\omega)$  and  $\mathcal{PR}_i(\omega')$  induce the same probability distribution on  $\Sigma_{-i} \times T_{-i}$ .

A BFK-like structure  $M$  whose state space is isomorphic to  $\vec{\Sigma} \times \vec{T}$  is *complete* if, for every distribution  $\mu_i$  over  $\Sigma_{-i} \times T_{-i}$ , (where the measurable sets are the ones induced by the isomorphism  $h$  and the measurable sets  $\mathcal{F}$  on  $\Omega$ ), there is a state  $\omega$  in  $M$  such that the probability distribution on  $\Sigma_{-i} \times T_{-i}$  induced by  $\mathcal{PR}_i(\omega)$  is  $\mu_i$ .

**Proposition 5.3:**  $M^c$  is a complete BFK-like structure.

We now would like to show that every measurable complete BFK-like structure is the canonical model. This is not quite true because states can be duplicated in an interactive structure. This suggests that we should try to show that the canonical structure can

be embedded in every measurable complete structure. We can essentially show this, except that we need to restrict to *strongly measurable* complete structures, where a structure is strongly measurable if it is measurable and the only measurable sets are those defined by  $\mathcal{L}_4$  formulas (or, equivalently, the set of measurable sets is the smallest set that contains the sets defined by  $\mathcal{L}_3$  formulas). (We explain the need for strong measurability in the full paper.)

**Theorem 5.4:** *If  $M$  is a strongly measurable complete BFK-like structure, then the canonical structure can be embedded in  $M$ .*

## 6 Discussion

We have provided a logical formula that captures the intuition that “all a player knows” is that players satisfy appropriate rationality assumptions. Our formalization of “all player  $i$  knows is  $\varphi$ ” is in terms of a language  $\mathcal{L}$ : roughly speaking, we require that  $i$  assigns positive probability to all formulas  $\psi \in \mathcal{L}$  that are consistent with  $\varphi$ . We showed that when  $\mathcal{L}$  expresses statements about the strategies played by the other players, our logical formula characterizes strategies that are iterated admissible (i.e., survive iterated deletion of weakly dominated strategies). On the other hand, when  $\mathcal{L}$  is less expressive (e.g., empty) the same logical formula instead characterizes strategies surviving iterated deletion of strictly dominated strategies (and thus also characterizes rationalizable strategies). Thus, the expressiveness of the language used by the players to describe their beliefs about other players can be viewed as affecting how the game is played. We plan to consider the effect of the players using other languages to describe their beliefs. For example, we are interested in the solution concept that arises if the language includes  $RAT_i$  for each player  $i$  but does not include  $play_i(\sigma)$ , so that players can talk about the rationality of other players without talking about the strategies they use.

We would also like to consider other ways of incorporating restrictions on how players form beliefs about other players. For example, we could restrict players’ beliefs to be consistent with a theory  $\mathcal{T}$ ; namely, we might require  $i$  to assign positive probability to all and only formulas  $\psi$  consistent with both rationality and  $\mathcal{T}$  (or, essentially equivalently, to all and only formulas that can be satisfied in some class of Kripke structures).

Such restrictions provide a straightforward way of capturing players’ prior beliefs about other players. They may also be used to capture the way “boundedly rational” players reason about each other. But what are “natural” restrictions? And how do such restriction affect how the game will be played? We leave an exploration of these questions for future research.

## Acknowledgements

The first author is supported in part by NSF grants ITR-0325453, IIS-0534064, and IIS-0812045, and by AFOSR grants FA9550-08-1-0438 and FA9550-05-1-0055. The second author is supported in part by NSF CAREER Award CCF-0746990, AFOSR Award FA9550-08-1-0197, BSF Grant 2006317 and I3P grant 2006CS-001-0000001-02. The second author wishes to thank Silvio Micali for helpful discussions and his excitement about the research direction.

## References

- Brandenburger, A., A. Friedenberg, and J. Keisler (2008). Admissibility in games. *Econometrica* 76(2), 307–352.
- Halpern, J. Y. and G. Lakemeyer (2001). Multi-agent only knowing. *Journal of Logic and Computation* 11(1), 41–70.
- Levesque, H. J. (1990). All I know: a study in autoepistemic logic. *Artificial Intelligence* 42(3), 263–309.
- Mas-Colell, A., M. Whinston, and J. Green (1995). *Microeconomic Theory*. Oxford, U.K.: Oxford University Press.
- Osborne, M. J. and A. Rubinstein (1994). *A Course in Game Theory*. Cambridge, Mass.: MIT Press.
- Pearce, D. G. (1984). Rationalizable strategic behavior and the problem of perfection. *Econometrica* 52(4), 1029–1050.
- Rajan, U. (1998). Trembles in the Bayesian foundation of solution concepts. *Journal of Economic Theory* 82, 248–266.
- Tan, T. and S. Werlang (1988). The Bayesian foundation of solution concepts of games. *Journal of Economic Theory* 45(45), 370–391.