

Respect for Public Preferences and Iterated Backward Inference

Oliver Schulte
School of Computing Science
Department of Philosophy
Simon Fraser University
Vancouver, Canada
oschulte@cs.sfu.ca

June 3, 2003
Draft - Please don't quote

Abstract

An important approach to game theory is to examine the consequences of beliefs that rational agents may have about each other. This paper considers *respect for public preferences*. Consider an agent A who believes that B strictly prefers an option a to an option b . Then A *respects* B 's preference if A considers the choice of a “infinitely more likely” than the choice of B ; equivalently, if A assigns probability 1 to the choice of a given that B chooses a or b . Respect for *public* preferences requires that if it is common belief that B prefers a to b , then it is common belief that all other agents respect that preference. Along the lines of Blume, Brandenburger and Dekel [4] and Asheim [1], I treat respect for public preferences as a constraint on lexicographic probability systems. The main result is that if respect for public preferences and perfect recall obtains, then players choose in accordance with Iterated Backward Inference. Iterated Backward Inference is a procedure that generalizes standard backward induction reasoning for games of both perfect and imperfect information. From Asheim's characterization of proper rationalizability [1] it follows that properly rationalizable strategies are consistent with respect for public preferences; hence strategies eliminated by Iterated Backward Inference are not properly rationalizable.

1 Introduction and Overview

Game theory provides a general formalism for representing social and economic interactions. The question arises what we can predict about the behaviour of the agents in a given situation. A *solution concept* gives a formal answer to this question, by associating a set of game plays—the “solution set”—with

a given game matrix or game tree. An important line of research examines epistemic assumptions that validate a given solution concept (cf. [19]). What conditions imply that the predictions of the solution concept are correct, in the sense that players will choose in accordance with the solution concept? For example, in their seminal work on rationalizability, Pearce and Bernheim examined the assumption that it is common knowledge among the agents that they each maximize subjective expected utility [14], [2]. Their work showed that iteratively eliminating strictly dominated strategies is a procedure that derives the predictions of this assumption exactly, in the sense that an outcome of the game is consistent with common knowledge of expected utility maximization if and only if the outcome is not eliminated.

This paper examines the implications of *respect for public preferences*. Blume *et al.* introduced the concept of respect for preferences [4, Def. 4] to characterize Myerson’s “proper equilibrium” [12]. Let us assume that each agent’s choices are based on a lexicographic probability system (LPS) $\rho = (\rho^1, \dots, \rho^k)$, where each ρ^i is a probability measure over the choices of other agents. According to Blume, Brandenburger, and Dekel, “the first component of the LPS [i.e., ρ^1] can be thought of as representing the player’s primary theory...” [4, page 82]. In a lexicographic probability system, the probability of an event E may be defined even conditional on an event E' that the agent believes not to obtain. If ρ_B is the LPS of agent B , then for any choice of an agent A between two options a and b , we may consider the conditional LPS $\rho_B|_{\{a, b\}}$. According to Blume *et al.*, B respects the preferences of A if $[\rho_B^1|_{\{a, b\}}](a) = 1$ whenever A strictly prefers option a to b . Intuitively, the “primary hypothesis” of B is that A chooses a , given that A chooses either a or b and prefers a . For example, suppose that agent A has three options, \$300, \$200, \$100. Then if A prefers \$200 to \$100, respect for preferences requires that $[\rho_B^1|_{\{\$200, \$100\}}](\$200) = 1$.

Asheim [1] introduced a weaker condition: according to his definition, B respects the preferences of A if $[\rho_B^1|_{\{a, b\}}](a) = 1$ whenever B believes (with certainty - see Section 5.2) that A strictly prefers option a to b . We may think of the definition of Blume *et al.* as a special case where B ’s beliefs about A ’s preferences are true—as they well may be at equilibrium.

Respect for *public* preferences requires common belief that $[\rho_B^1|_{\{a, b\}}](a) = 1$ whenever it is *common belief* among all agents that A strictly prefers option a to b . An event is common belief among the agents if all agents believe that it obtains, all agents believe all agents believe that it obtains, etc. Preferences that are common belief are “public” in the sense that all agents are aware of them. Assuming that agents know their own preferences, then if A believes that she prefers a to b , this is indeed the case, and hence public preferences are true preferences.

This paper is a formal investigation of what respect for public preferences implies about agents’ behaviour. More specifically, I derive consequences of common belief in the following assumptions:

Respect for Public Preferences If it is common belief that player A strictly prefers option a over b , then it is common belief that $[\rho_B^1|_{\{a, b\}}](a) = 1$

for each player $B \neq A$.

Full Lexicographic Rationality It is common belief that each player maximizes lexicographic expected utility, with respect to an LPS with *full support*. A lexicographic probability system ρ has full support if every nonempty event receives positive probability at some measure in ρ .

For short, in this introduction I refer to these two assumptions as *Respect for Public Preferences*. I specify an iterated elimination procedure that computes consequences of Respect for Public Preferences, which I term Iterated Backward Inference (IBI). In two special cases, IBI coincides with other well-known algorithms. First, in a game of perfect information with a unique backward induction solution, the result of IBI is that solution. Second, suppose we represent a strategic form game as a game tree with two information sets, one for each player with moves corresponding to strategies. Then IBI amounts to the following procedure: First, eliminate all weakly dominated strategy. Then iteratively eliminate all strictly dominated strategies. This procedure is known as the *Dekel-Fudenberg* procedure [6]. Thus Iterated Backward Inference generalizes at once both standard backward induction and the Dekel-Fudenberg procedure.

The argument for these results, and more generally for the validity of IBI for computing consequences of Respect for Public Preferences, rests on two key facts. First, *lexicographic rationality implies Sequential Admissibility*. I define the notion “admissibility at an information set” precisely in Section 10 below. It turns out that sequential admissibility follows from lexicographic rationality in the sense that if a strategy is inadmissible at an information set, then it does not maximize lexicographic expected utility. The only assumption required for this result is that each player has perfect recall.

Second, *lexicographic beliefs satisfy Backward Inference*. Loosely speaking, the principle of backward inference is that agents “look ahead and reason back”. More precisely, I consider the following principle: Think of players as following a strategy at various information sets. Let I, I_B be two information sets such that I_B belongs to player B and I_B follows I ; that is, there is a play sequence on which I is reached before I_B . Then if player A “looks ahead” to I_B and believes that player B is unlikely to follow a certain strategy s at information set I_B , then player A “reasons back” and believes that player B is unlikely to follow the strategy s at information set I . In terms of lexicographic probability systems, we have that if $[\rho_A^1|I_B](s) = 0$, then $[\rho_A^1|I](s) = 0$. I show that lexicographic probability systems satisfy the backward inference principle, provided that each player has perfect recall. Thus rather than being a separate, additional principle, backward inference follows from the structure of lexicographic beliefs.

Applying Asheim’s characterization of Schumacher’s concept of proper rationalizability [1], [16] it is easy to show that properly rationalizable strategies are consistent with Respect for Public Preferences. Hence if IBI eliminates a strategy s_i , then s_i is not properly rationalizable. So IBI can be used to find strategies that are not properly rationalizable.

The paper is organized as follows. Sections 2 and 3 define standard game-theoretic notions such as game trees and strategies, and review definitions pertaining to lexicographic probability systems. Sections 4 and 5 formalize a number of epistemic assumptions, particularly Respect for Public Preferences, and discuss the relationship between this principle and proper rationalizability. The remainder of the paper investigates the consequences of these assumptions. Sections 6 and 7 define Iterated Backward Inference, establish its correctness (i.e., soundness) and show existence for finite games—that is, in finite games some strategy profile is guaranteed to survive Iterated Backward Inference. Section 8 demonstrates that in games of perfect information in which backward induction provides a unique solution, Iterated Backward Inference agrees with backward induction. Next I prove a number of auxiliary results required for establishing the soundness of Iterated Backward Inference. Perfect recall is crucial for reasoning from an extensive form of a game to its normal form; Section 9 establishes a number of basic properties of game trees with perfect recall. Section 10 shows that under perfect recall, lexicographic rationality entails sequential admissibility.

Unless otherwise stated, proofs appear in Section 13.

2 Preliminaries: Game Trees and Strategies

This section gives the standard definition of sequential games; I also use the terms *extensive form games* or *game trees*. I employ the formulation from [13]. A key notion in what follows is that of a sequence. Almost all the sequences I consider in this paper are sequences of actions. I denote actions throughout by the letter a and variants. I write $x = a_1, \dots, a_n$ to indicate the finite sequences whose i -th member is a_i , and similarly I write $h = a_1, \dots, a_n, \dots$ for infinite sequences. If $x = a_1, \dots, a_n$ is a finite sequence of n actions, the concatenation $x * a = a_1, \dots, a_n, a$ yields a finite sequence of length $n + 1$. The relation “ x is a prefix of x' ” is a partial ordering of sequences. I write $x \leq x'$ to indicate that sequence x is a prefix of x' , and $x' < x$ to indicate that $x \leq x'$ and $x \neq x'$. The notation $x' \geq x$ and $x' > x$ denotes the dual notions (x is a prefix of x' , or x' extends x).

Now we are ready to define a sequential game.

Definition 1 A *sequential game* T is a tuple $\langle N, H, \text{player}, \{I_i\}, \{u_i\} \rangle$ whose components are as follows.

1. A finite set N (the set of **players**).
2. A set of H of sequences (**histories**) that satisfies the following three properties.
 - (a) The empty sequence \emptyset is a member of H .
 - (b) If x is in H , then every initial segment of x is in H .

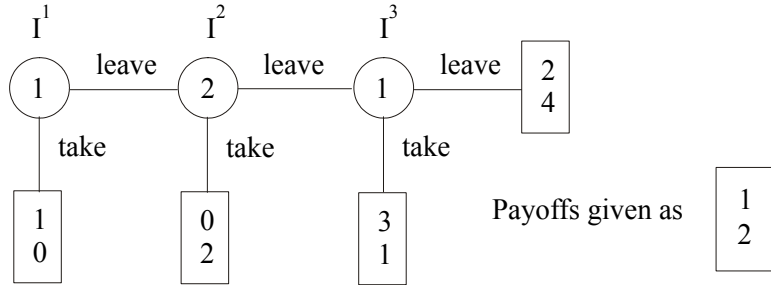


Figure 1: The 3-move version of the Centipede game with players A and B .

(c) If h is an infinite sequence such that every finite initial segment of h is in H , then h is in H .

Each component of a history in H is an **action** taken by a player. A finite member of H is a **node**. A node $x \in H$ is **terminal** if there is no history $x' \in H$ such that $x' > x$. All infinite histories are terminal as well. The set of terminal histories is denoted by Z . The set of actions available at a node x is denoted by $A(x) = \{a : x*a \in H\}$.

3. A function player that assigns to each nonterminal node a member of N . The function player determines which player takes an action at node x .
4. For each player $i \in N$ an **information partition** \mathcal{I}_i defined on $\{x \in H : \text{player}(x) = i\}$. An element I_i of \mathcal{I}_i is called an **information set** of player i . We require that if x, x' are members of the same information set I_i , then $A(x) = A(x')$. For each information set I_i , I let $A(I_i)$ be the set of actions available at the nodes in I_i . I let $A_i =_{df} \cup \{A(x) : \text{player}(x) = i\}$ denote the set of actions of player i .
5. For each player $i \in N$ a **payoff function** $u_i : Z \rightarrow R$ that assigns a real number to each terminal node.

Though the general notion of an extensive form game permits “chance” moves by “nature”, Definition 1 does not include chance moves. Another simplification that I make throughout the paper is to consider only 2-player games, that is, I take $N = \{1, 2\}$. It is straightforward to generalize the results in this paper to games with chance moves and any finite number of players, but doing so gives rise to technical complications that do not illuminate the main issues.

I illustrate Definition 1 in two extensive form games, the three-move version of the well-known Centipede Game (Figure 1) and a game due to Kohlberg (Figure 2) [10]. In my notation, the 7 nodes in the Centipede Game are $\emptyset, \text{take}, \text{leave}, \text{leave} * \text{take}, \text{leave} * \text{leave}, \text{leave} * \text{leave} * \text{take}, \text{leave} * \text{leave} * \text{leave}$. The terminal nodes comprise $\text{leave} * \text{leave} * \text{leave}$ and all sequences ending in take . Each of the three information sets I^1, I^2, I^3 is a singleton, which

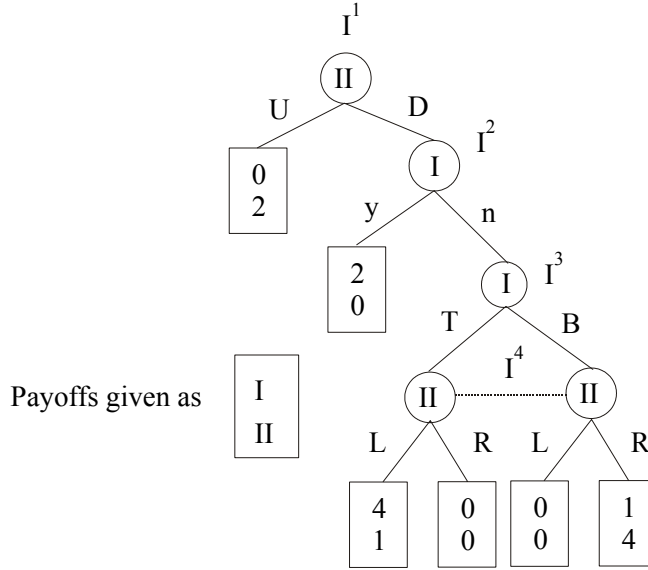


Figure 2: An extensive form game due to Kohlberg

makes the Centipede game a game of **perfect information**; for example, $I^2 = \{leave * leave\}$. Kohlberg's game is a game of imperfect information because we have that $I^4 = \{D * n * T, D * n * B\}$.

If i denotes a player, I write $-i$ for the opponent of player i ; so $-1 = 2$ and $-2 = 1$ when 1, 2 refer to players. A **strategy** for player i is a function $s_i : \{x \in H : player(x) = i\} \rightarrow A_i$, such that (1) $s_i(x) \in A(x)$ for all nodes x belonging to player i , and (2) if $I(x) = I(x')$, then $s_i(x) = s_i(x')$. I write $S_i(T)$ for the **set of strategies** of player i in T . A strategy pair (s_1, s_2) of players 1 and 2 respectively determines a unique terminal history that I denote by $play(s_1, s_2)$. I extend the utility functions u_i to strategy pairs by defining $u_i(s_1, s_2) =_{df} u_i(play(s_1, s_2))$.

I assume throughout the paper that *every node is reachable by some pair of strategies*. That is, I assume that for all nodes $x \in H$, there is a strategy pair (s_1, s_2) such that x is reached along $play(s_1, s_2)$ (cf. [9, Sec.2.1]).

To illustrate, in the Centipede Game there are two strategies for player 2, which I denote by l_2 and t_2 where $l_2(\emptyset * leave) = leave$, and $t_2(\emptyset * leave) = take$; thus if T denotes the Centipede Game, then $S_2(T) = \{l_2, t_2\}$. Player 1 has four strategies, specifying choices at I^1 and I^3 . We may use the tuples in $\{l, t\} \times \{l, t\}$ to denote these strategies, such that for example $lt(\emptyset) = leave$, and $lt(\emptyset * leave * leave) = take$. The play sequence resulting from the strategy pair (lt, t_2) is $\emptyset * leave * leave * take$; in my notation $play(lt, t_2) = \emptyset * leave * leave * take$. Thus $u_1(lt, t_2) = 0$, and $u_2(lt, t_2) = 2$.

A matrix whose rows correspond to strategies for player 1 and columns to

1/2	l_2	t_2
tt	1,0	1,0
tl	1,0	1,0
lt	3,1	0,2
ll	2,4	0,2

Table 1: The strategic form of the Centipede Game

I/II	UL	UR	DL	DR
yT	0,2	0,2	2,0	2,0
yB	0,2	0,2	2,0	2,0
nT	0,2	0,2	4,1	0,0
nB	0,2	0,2	0,0	1,4

Table 2: The strategic form of Kohlberg's game

strategies for player 2 gives the **normal form**, or **strategic form**, of a game tree. Table 1 shows the strategic form of the Centipede game. In Kohlberg's game, each player has four strategies. Table 2 shows the normal form of Kohlberg's game.

A crucial question in the developments below is how player i ranks the strategies of her opponent $-i$ given the information in some information set I . To describe this formally, I map information sets into sets of strategies as follows. First, define $[I] = \{(s_1, s_2) : \text{play}(s_1, s_2) \text{ intersects } I\}$; thus $[I]$ denotes the set of strategy pairs that are consistent with the information I . With respect to player i 's uncertainty space S_{-i} , the information in an information set I corresponds to the set of strategies $[I]_{-i}$ of player $-i$ that are consistent with I ; formally I define $[I]_{-i} = \{s_{-i} : \exists s_i.(s_i, s_{-i}) \in [I]\}$. It is useful to consider, for a strategy s_i of player i , the set of strategies of the other player that are consistent with s_i at a given information set. Accordingly, I define $\text{cons}(s_i, I) = \{s_{-i} : (s_i, s_{-i}) \in [I]\}$. To illustrate, in the Centipede game of Figure 1, $[I^2] = \{(ll, l_2), (lt, l_2), (ll, t_2), (lt, t_2)\}$, so $[I^2]_1 = \{ll, lt\}$ and $[I^2]_2 = \{l_2, t_2\}$. For information set I^3 we have that $\text{cons}(tt, I^3) = \emptyset$, and $\text{cons}(lt, I^3) = \{l_2\}$, while for player 2 $\text{cons}(t_2, I^3) = \emptyset$ and $\text{cons}(l_2, I^3) = \{ll, lt\}$.

Following Kaneko and Kline [9, Sec.2.1], I define $I > I'$ iff there is $x \in I, x' \in I'$ such that $x > x'$, and $I \geq I'$ by $I > I'$ or $I = I'$. Thus $I > I'$ holds just in some play sequence reaches I' and then I . I also employ the dual notions $I < I'$ and $I \leq I'$. For example, in the Kohlberg game of Figure 2, we have that $I^j > I^4$ for all information sets $I^j \neq I^4$.

ρ_B^3/X	\$100	\$200	\$300
ρ_B^3	1	0	0
ρ_B^2	0	1	0
ρ_B^1	0	0	1

Table 3: A lexicographic probability system ρ_B of length 3, with states of the world $X = \$100, \$200, \$300$

ρ_2^3/X	tt	tl	lt	ll
ρ_2^3	0	0	1/2	1/2
ρ_2^2	0	0	1	0
ρ_2^1	1/2	1/2	0	0

Table 4: A lexicographic probability system ρ_2 representing the beliefs of player 2 about the strategies of player 1 in the Centipede Game

3 Preliminaries: Lexicographic Expected Utility

Let X be a finite set of points. A **lexicographic probability system** over X is a finite sequence or vector $\rho = (\rho^1, \rho^2, \dots, \rho^k)$, where each ρ^j is a probability measure over X . As indicated, for a given LPS ρ I write ρ^j for the j -th probability measure in the sequence ρ . I let $|\rho|$ denote the **length** of ρ . I also write $\rho(j)$ for the **support** of ρ^j , that is $\rho(j) = \{x \in X : \rho^j(x) > 0\}$. An LPS ρ has **full support** iff $\cup\{\rho(i) : 1 \leq i \leq |\rho|\} = X$. Thus if ρ has full support, then every point x is in the support of some probability measure in ρ .

Following Blume *et al.* [4, Definition 2], I write $x \geq_\rho x'$ if $\min\{j : x \in \rho(j)\} \leq \min\{j : x' \in \rho(j)\}$, and $x >_\rho x'$ if the inequality is strict. Informally, $x >_\rho x'$ means that the agent considers x more plausible than x' , in that x is consistent with a “lower-order” belief than x' .

To illustrate these definitions, let $X = \{\$100, \$200, \$300\}$ as in the simple choice problem from Section 1. Then a probability measure p corresponds to a ternary vector (e.g., if $p(\$100) = 1$, then p corresponds to $(1, 0, 0)$), and we may have the LPS $\rho_B = [(0, 0, 1), (0, 1, 0), (0, 0, 1)]$ representing the beliefs of agent B —see Table 3. Then $\$300 >_{\rho_B} \$200 >_{\rho_B} \$100$. Table 4 shows an LPS ρ_2 that might represent player 2’s beliefs about 1’s strategies in the Centipede Game. Here $tt >_{\rho_2} lt >_{\rho_2} ll$. Both lexicographic systems ρ_B and ρ_2 have full support.

Let S be a set of acts, and let $u : S \times X \rightarrow R$ be a utility function. The expected utility of an act s with respect to probability p and utility u is denoted by $EU(s, p, u)$ and defined to be $EU(s, p, u) =_{df} \sum_{x \in X} p(x) \times u(s, x)$. The **lexicographic expected utility** of an act s with respect to an LPS ρ is a vector of real numbers (expected utilities) defined as $LEU(s, \rho, u) =_{df} (EU(s, \rho^1, u), EU(s, \rho^2, u), \dots, EU(s, \rho^k, u))$ where $k = |\rho|$. For two vectors \mathbf{u}, \mathbf{u}' of real numbers, I let $\mathbf{u} \geq \mathbf{u}'$ denote the lexicographic ordering of the two vectors. Then an act s **maximizes lexicographic expected utility**

ρ_2^i/X	tt	tl	lt	ll	$EU(l_2, \rho_2, u_2)$	$EU(t_2, \rho_2, u_2)$
ρ_2^3	0	0	1/2	1/2	2.5	2
ρ_2^2	0	0	1	0	1	2
ρ_2^1	1/2	1/2	0	0	0	0

Table 5: The lexicographic expected utility of player 2’s strategies in the Centipede game with utility function u_2 for player 2, given LPS ρ_2 from Table 4

given ρ, u if $LEU(s, \rho, u) \geq LEU(s', \rho, u)$ for all $s' \in S$. A preference ordering \succeq over S is **represented** by a pair (ρ, u) if for all options s, s' we have that $s \succeq s' \iff LEU(s, \rho, u) \geq LEU(s', \rho, u)$. I say that an agent with preference ordering \succeq is a maximizer of lexicographic expected utility if \succeq is represented by some pair (ρ, u) . To illustrate, Table 5 shows the lexicographic expected utility of the strategies l_2 and t_2 for Player 2 in the Centipede game, given ρ_2 . In this example, $LEU(l_2, \rho_2, u_2) = (0, 1, 2.5)$, and $EU(t_2, \rho_2, u_2) = (0, 2, 2)$, so if (ρ_2, u_2) represent the preferences of Player 2, then $t_2 \succ l_2$ because $EU(t_2, \rho_2^2, u_2) = 2 > 1 = EU(l_2, \rho_2^2, u_2)$.

Many decision theorists [3] view it as an attractive feature of lexicographic probability systems with full support that they enforce the principle of *weak dominance*: if an act s weakly dominates another act s' , then the lexicographic expected utility of s is greater than the lexicographic expected utility of s' . For future reference, I record this fact as a formal lemma. Given a utility function $u : S \times X \rightarrow R$, I say that an act s **weakly dominates** an act s' if (1) for all points x , we have $u(s, x) \geq u(s', x)$, and for some point x' , we have $u(s, x') > u(s', x')$. For example, in the Centipede game, lt weakly dominates ll (see Table 1). The following lemma is proven in [3, Theorem 4.2].

Lemma 2 *Let ρ be a lexicographic probability system over X with full support. Suppose that act s weakly dominates s' with respect to utility function $u : S \times X \rightarrow R$. Then $LEU(s, \rho, u) > LEU(s', \rho, u)$.*

3.1 Conditional Lexicographic Beliefs and Preferences

I next define the conditional LPS $\rho|P$ given a nonempty event P . As usual, if p is a probability measure over X , and $P \subseteq X$ an event such that $support(p) \cap P \neq \emptyset$, the conditional probability $p|P$ is defined by $[p|P](x) = p(x)/p(P)$ if $x \in P$, and $[p|P](x) = 0$ otherwise. Intuitively, to obtain the conditionalized probability system $\rho|P$, we first delete all probability measures ρ^j such that $support(\rho^j) \cap P = \emptyset$, and condition the remaining probability measures on P . For example, conditioning ρ_2 on $\{lt, ll\}$ yields $\rho_2|\{lt, ll\}$ displayed in Table 6. The formal definition of $\rho|P$ is as follows.

1. $o(1) = \min\{1 \leq k \leq |\rho| : \rho(k) \cap S \neq \emptyset\}$. If ρ has full support, $o(1)$ is well-defined. And $(\rho|P)^1 = (\rho^{o(1)})|P$.

$[\rho_2^o \{lt, ll\}]/X$	tt	tl	lt	ll
$[\rho_2^2 \{lt, ll\}]$	0	0	1/2	1/2
$[\rho_2^1 \{lt, ll\}]$	0	0	1	0

Table 6: The result of conditioning ρ_2 in Table 4 on $\{lt, ll\}$

2. $o(n+1) = \min\{o(n) < k \leq |\rho| : \rho(k) \cap S \neq \emptyset\}$. If there is no such n , then $|\rho|P| = n$. Otherwise $(\rho|P)^{n+1} = (\rho^{o(n+1)})|P$.

If $\rho|P$ does not have full support, the conditional probability $\rho|P$ may not be well-defined. According to decision theorists, another attractive feature of lexicographic probability systems with full support is that conditional probabilities are well-defined for any event [3]. We can use conditional probabilities to represent conditional preferences.

Definition 3 *Let X be a set of states of the world and O a set of options. Suppose that (ρ, u) represent a preference ordering \succeq over O , where ρ has full support. Then for each nonempty event $P \subseteq X$, the conditional preference $\succeq^{|P}$ is defined by: $a \succeq^{|P} b \iff LEU(a, \rho|P, u) \geq LEU(b, \rho|P, u)$.*

For an event $P \subseteq X$, I write \bar{P} for the complement of P in X . I remark that in the paper of Blume, Brandenburger and Dekel [3], conditional preference is defined differently, namely in the standard decision-theoretic manner: to compare a and b conditional on event P , compare two acts a', b' that yield the same outcome in all states of the world in \bar{P} , and that agree with a and b respectively on states of the world in P . Blume, Brandenburger and Dekel then prove that for an LPS with full support, conditional lexicographic expected utility characterizes conditional preferences as indicated in Definition 3 [3, Th.4.3].

I will make use of the fact that lexicographic preferences satisfy the “sure thing” principle: improving the expected payoff of an option a over a range of possibilities P without affecting the expected payoff over \bar{P} yields an option a' that is preferred to a overall. Formally, we have the following result.

Proposition 4 (Blume, Brandenburger and Dekel) *Suppose that (ρ, u) represent a preference ordering \succeq , where ρ has full support, and let P be a nonempty event. For all options a, b , if $a \succ^{|P} b$ and $a \sim^{\bar{P}} b$, then $a \succ b$.*

The proposition is an immediate consequence of Theorem 4.1 (i) in [3].

3.2 Revision of Lexicographic Beliefs

I introduce a new operation $\rho * P = [\rho|P](1)$, which assigns to each event P the support of the first probability measure in $\rho|P$. I refer to this operation as the **revision** of ρ on P . For example, in the LPS ρ_2 above, we have that $\rho_2 * \{lt, ll\} = \{lt\}$ (see Table 6). A revision on P can be thought of as representing the “primary theory” or “first-order beliefs” of the agent given the information

P . The following simple lemma provides an intuitive interpretation of revision: the revision of ρ on P selects exactly the states of the world x in P that are maximal in the “plausibility” ordering \geq_ρ . I omit the straightforward proof.

Lemma 5 *Let ρ be an LPS over X , and let $P \subseteq X$. Then for all points $x \in P$, it is the case that $x \in (\rho * P) \iff \forall x' \in P. x \geq_\rho x'$.*

As Stalnaker has noted [20, fn.12], lexicographic probability systems are closely related to belief revision structures that feature in the well-known AGM belief revision theory [7], [17]. Each LPS induces a revision operator $+$ defined by $\rho(1) + P =_{df} \rho * P$. If we interpret this operation as a revision of the agent’s “primary theory” $\rho(1)$ on the information P , it is easy to verify that the revision satisfies the well-known AGM axioms for minimal belief change. A difference is that in the AGM theory, a belief revision operator is a *binary* function from “current beliefs” K and “new information” P to new beliefs; in contrast, the revision associated with an LPS is a *unary* function of information P . (Hans Rott discusses advantages and disadvantages of unary vs. binary belief revision operators [15].)

4 General Epistemic Assumptions

Consider a 2-player game $G = \langle S_1, S_2, u_1, u_2 \rangle$, with sets of options S_1 and S_2 , and utility functions u_1, u_2 defined on $S_1 \times S_2$. Let W be a set of states of the world. A given state of the world w associates the following elements with each player i :

1. a strategy choice $choice_i^w \in S_i$
2. a preference ordering \succeq_i^w over the options S_i
3. a LPS ρ_i^w over S_{-i} , and hence a weak ordering $\geq_{\rho_i^w}$ over S_{-i} ; I write more concisely $\geq_i^w =_{df} \geq_{\rho_i^w}$.
4. a belief operator B_i^w . If A is an assertion about the game G , then $B_i(A)$ expresses the fact that in w , player i believes A .

One may take the belief operator B_i as given or interpret it in various ways, for example such that $B_i(A)$ represents probability 1 belief in A [18], or that $B_i(A)$ is the “first-order belief” of an agent in a lexicographic probability system, for example an LPS over a type space [11], [5], or that $B_i(A)$ corresponds to “certain belief” [1] (\bar{A} receives probability 0 in every probability measure of an LPS of player i ; see Section 5.2.) For both simplicity and generality, I do not interpret the belief operator further, but instead state axiomatically my assumptions about it. The theorems in this paper hold for any concept of belief that satisfies my axioms.

In what follows, I consider the implications of various conditions on the epistemic elements listed above. Using the techniques of epistemic logic, it would

be possible to develop a precise formal language for stating such assumptions [8]. For most of this paper, a semi-formal understanding is sufficient. I begin with a number of basic epistemic principles. I will indicate throughout the paper which of my results depends on which of these principles.

Definition 6 *Basic Epistemic Principles*

1. (*Lexicographic rationality*) ρ_i, u_i represent \succeq_i .
2. (*Full Support*) ρ_i has full support.
3. (*Preference Maximization*) If $\text{choice}_i = s_i$, then $\forall s'_i, s_i \succeq_i s'_i$.
4. (*Preference Introspection*) If $B_i(s_i \succeq s'_i)$, then $s_i \succeq s'_i$.

I use the notion of **common belief**. If A is an assertion about the game G , I write $CB(A)$ to denote that A is common belief among the players. Define the set of belief assertions about A , denoted by $F(A)$, as follows.

1. $B_i(A)$ and $B_{-i}(A)$ are in F .
2. If F' is in $F(A)$, then $B_i(F')$ and $B_{-i}(F')$ are in $F(A)$.
3. No other expression is in $F(A)$.

Then $CB(A)$ holds in a state of the world iff all belief assertions $F(A)$ about A hold. Note that it follows from the definition that if $CB(A)$ holds, then so does $CB(B_i(A))$ for each player i . I make the following assumptions about common belief throughout this paper.

1. Common belief is *closed under implication*, that is, if $CB(\text{if } A \text{ then } B)$ and $CB(A)$ hold, then $CB(B)$ holds as well.
2. Mathematical and logical truths are common belief.
3. All aspects of the game, or game tree for extensive form games, are common belief among the players. In particular, the utility functions are common belief, and thus there is no uncertainty about payoffs.

5 Respect for Preferences

This section formalizes my key assumption, respect for public preferences, and clarifies its relationship with previous work, particularly with Asheim's notion of respect for preferences and proper consistency.

5.1 Respect for Public Preferences

Let us return to the example of Section 1 in which an agent A has three options, \$300, \$200, \$100. Suppose that an agent B believes that A 's preference ranking is $\$300 \succ \$200 \succ \$100$. Then if B respects preferences, the first-order probability measure of B 's lexicographic belief system conditional on the event $\{\$200, \$100\}$ assigns probability 1 to \$200. In terms of the revision of lexicographic beliefs described in Section 3, the revision of B 's beliefs on the event $\{\$200, \$100\}$ entails that A chooses \$100. As a general principle, we have the following axiom.

Axiom 7 *For each agent i , if $B_i(s_{-i} \succ s'_{-i})$, then $\rho_i * (\{s_{-i}, s'_{-i}\}) = \{s_{-i}\}$.*

Axiom 7 constrains revisions on binary choices of another agent, which in the game-theoretic setting are revisions on a set containing two states of the world; it is easy to see that the $>_\rho$ ordering associated with an LPS ρ determines the result of such revisions. Hence we have the following lemma.

Lemma 8 *Let ρ be a lexicographic probability system ρ , and let x, x' be any two states of the world. Then $\rho * \{x, x'\} = \{x\} \iff x >_\rho x'$.*

If $B_i(s_{-i} \succ s'_{-i})$, then by Axiom 7, $\rho_i * (\{s_{-i}, s'_{-i}\}) = \{s_{-i}\}$. So by Lemma 8, $B_i(s_{-i} \succ s'_{-i})$ implies that $s_{-i} >_i s'_{-i}$. Thus if agent i believes that agent $-i$ prefers option s_{-i} to s'_{-i} , then agent i ranks s_{-i} higher than s'_{-i} , or “considers s_{-i} infinitely more likely than s'_{-i} ”. A weaker assumption is that this principle holds only for preferences that are *public*, in the sense that they are common belief among the agents. My next axiom asserts that for public preferences that are common belief, it is also common belief that the preferred option is ranked above the dispreferred one.

Axiom 9 (*Respect for Public Preferences*) *For each agent i , if $CB(s_i \succ_i s'_i)$, then $CB(s_i >_{-i} s'_i)$.*

Axiom 9 is the crucial assumption for deriving the validity of the Iterated Backward Inference procedure that I describe in Section 6; I refer to Axom 9 as Respect for Public Preferences. Given our assumptions about common belief, it is possible to derive Axiom 9 from common belief in Axiom 7.

Lemma 10 *Common belief in Axiom 7 implies Respect for Public Preferences. In other words, common belief in Axiom 7 implies Axiom 9.*

The remainder of this paper investigates the consequences of Respect for Public Preferences, combined with the general epistemic assumptions laid out in Section 4. Before I start this investigation, I clarify the relationship between my epistemic assumptions and previous work, particularly Schuhmacher and Asheim's results on proper rationalizability. Reading the next subsection can be omitted without loss of continuity.

5.2 Proper Rationalizability and Respect for Preferences

Axiom 7 is very closely related to Asheim’s definition of “respecting preferences” [1, Sec. 4.1]. I outline Asheim’s definition to clarify the precise relationship; for more details see [1]. The definition is given in a semantic framework with types. A state of the world is a tuple (s_1, s_2, t_1, t_2) where s_i is a pure strategy for player i and t_i is a *type* of player i from the set T_i of types of player i . In Asheim’s framework, there are only finitely many types for each player [1, Def.1]. For each type $t_i \in T_i$ there is a preference relation \succeq^{t_i} over the pure strategies of player i , and an LPS ρ^{t_i} over the points in $S_{-i} \times T_{-i}$. For a given lexicographic system ρ over points X , define *certain belief* B_ρ by $B_\rho(E)$ iff $\text{supp}(\rho) \cap E = \emptyset$. In other words, E is certain belief given ρ iff E receives probability 0 in every measure in ρ . The dual *possibility operator* P_ρ associated with an LPS ρ is given by $P_\rho(E) \iff \neg B_\rho(\bar{E})$. Finally, define the events $[t_i] = \{(s_i, t_i) : s_i \in S_i\}$ and $[s_i] = \{(s_i, t_i) : t_i \in T_i\}$, where the events are in the space $S_i \times T_i$.

Asheim introduces a *cautiousness* condition for players’ beliefs. The event that player i is cautious is defined by: $(s_1, s_2, z_1, z_2) \in \text{cau}_i \iff$ for all types $t_{-i} \in T_{-i}$, if $P_{\rho^{z_i}}([t_{-i}])$, then $P_{\rho^{z_i}}(\{(s_{-i}, t_{-i})\})$ for all $s_{-i} \in S_{-i}$. So player i is cautious if for each type t_{-i} that i considers epistemically possible, we have that i considers all strategies $s_{-i} \in S_{-i}$ (i.e., all pairs (s_{-i}, t_{-i})) epistemically possible.

In this notation, Asheim’s definition of the event that player i respects preferences corresponds to: $(s_1, s_2, z_1, z_2) \in \text{resp}_i \iff$ for all pairs $(s_{-i}, t_{-i}), (s'_{-i}, t_{-i}) \in S_{-i} \times T_{-i}$, if $P_{\rho^{z_i}}([t_{-i}])$ and $s_{-i} \succ^{t_{-i}} s'_{-i}$, then $(s_{-i}, t_{-i}) \succ_{\rho^{z_i}} (s'_{-i}, t_{-i})$. To illustrate, suppose that player 1 considers epistemically possible a type t_2 of player 2 (i.e., $P_{\rho^{t_1}}([t_2])$), and that type t_2 prefers strategy s_2 to s'_2 . Then player 1 ranks the pair (s_2, t_2) higher than (s'_2, t_2) in his LPS ρ^{t_1} . The intuition is very close to that behind Axiom 7: given that player 2’s prefers s_2 to s'_2 (i.e., given that his type is t_2), player 1 ranks s_2 higher than s'_2 .

To see how my results relate to Asheim’s axioms, interpret the belief operator B_i as certain belief $B_{\rho^{t_i}}$, and the lexicographic probability system ρ_i as the marginal $\rho_m^{t_i}$ of ρ^{t_i} , which is defined by $(\rho_m^{t_i})^j(s_{-i}) =_{df} \sum_{t_{-i} \in T_{-i}} (\rho^{t_i})^j(s_{-i})$. With this interpretation, *cautiousness and respect for preferences imply Axiom 7*. To verify the implication, let $[s_i \succ_i s'_i] = \{(r_i, t_i) : s_i \succ^{t_i} s'_i\}$ be the event in $S_i \times T_i$ corresponding to player i ’s preference for s_i over s'_i , for each player i . Assume cautiousness and respect for preference and suppose that $B_i(s_i \succ s'_i)$ holds in a state of the world (s_1, s_2, t_1, t_2) , that is, $B_{\rho^{t_i}}([s_i \succ_i s'_i])$ holds. Consider the revision $\rho_m^{t_i} * (\{s_{-i}, s'_{-i}\})$. Let $r_{-i} \in \rho_m^{t_i} * (\{s_{-i}, s'_{-i}\})$; then there is a type t_{-i} such that $(t_{-i}, r_{-i}) \in \rho^{t_i} * ([s_{-i}] \cup [s'_{-i}])$, and $r_{-i} = s_i$ or $r_{-i} = s'_{-i}$. Since player i believes with certainty that $s_{-i} \succ_{-i} s'_{-i}$ (i.e., $B_{\rho^{t_i}}([s_{-i} \succ_{-i} s'_{-i}])$), it follows that $s_{-i} \succ^{t_{-i}} s'_{-i}$. By cautiousness, both (s_{-i}, t_{-i}) and (s'_{-i}, t_{-i}) are in the support of some measure in ρ^{t_i} . Since $s_{-i} \succ^{t_{-i}} s'_{-i}$, respect for preferences implies that $(s_{-i}, t_{-i}) \succ_{\rho^{t_i}} (s'_{-i}, t_{-i})$. Hence by Lemma 5, $(s'_{-i}, t_{-i}) \notin \rho^{t_i} * ([s_{-i}] \cup [s'_{-i}])$, and so $r_{-i} = s_{-i}$. Thus for all $r_{-i} \in \rho_m^{t_i} * (\{s_{-i}, s'_{-i}\})$, we have that $r_{-i} = s_{-i}$, as required by Axiom 7.

From a strict logical point of view, the converse is false: It is possible for Axiom 7 to hold even given cautiousness and respect for preferences. For example, if player i does not certainly believe in any preference of player $-i$, then Axiom 7 holds vacuously even if player i does not respect preferences. However, conceptually there appears to be very little difference between the two principles. For player i to respect preferences basically requires that conditional on the information that player $-i$ has some (epistemically possible) type t_{-i} , and hence a preference ordering \succeq^{t_i} , the LPS of player i satisfies Axiom 7 with respect to \succeq^{t_i} . It seems therefore that in any state in which Axiom 7 is common belief, respect for preferences ought to be common belief as well.

It is easy to see that cautiousness implies my principle of Full Support (Definition 6): if ρ^{t_i} is cautious, then $\rho_m^{t_i}$ has full support over S_{-i} . It is interesting to note, however, that certain belief in cautiousness is much stronger than certain belief in Full Support. Indeed, certain belief in cautiousness is inconsistent with certain belief in Preference Maximization (Definition 6), for any pair (t_i, s_i) such that s_i is not maximally preferred by type t_i violates Preference Maximization. Hence certain belief in Preference Maximization implies that $-i$ considers (t_i, s_i) impossible (i.e., $\neg P_{\rho^{t_{-i}}}(\{t_i, s_i\})$), which violates cautiousness for any type t_i that is subjectively possible for $-i$ (i.e., $P_{\rho^{t_{-i}}}([t_i])$ holds). Intuitively, if all mistakes are subjectively possible, then it cannot be subjectively impossible that an agent makes no mistake. One may of course have certain belief in cautiousness together with a weaker form of belief in Preference Maximization, for example “first-order belief” in the sense that $(\rho^{t_{-i}})^1$ assigns probability 1 to Preference Maximization.

Respect for public preferences (Axiom 9) is both logically and conceptually weaker than respect for preferences because it applies only to preferences that are common belief among the agents. Thus even if, for example, player 1 has certain belief that player 2 prefers s_2 to s'_2 , but this fact is not common (certain) belief, then Axiom 9 does not require player 1 to respect the preference of s_2 over s'_2 (i.e., the Axiom allows that $s'_2 >_{\rho^{t_1}} s_2$). As it turns out, the weak condition of respect for public preferences is sufficient to validate Iterated Backward Inference, the elimination procedure presented in this paper.

Asheim refers to the conjunction of cautiousness and respect for preferences (plus knowledge of the game structure) as *proper consistency* [1, Sec. 4.1]. A strategy s_i is *properly rationalizable* if s_i maximizes preferences in a state of the world in which there is common (certain) belief of proper consistency [1, Def. 2]. Asheim proves that this definition of proper rationalizability is equivalent to the previous one by Schuhmacher [1, Prop. 3], [16]. Since common belief in proper consistency entails common belief in Axiom 7, which in turn entails respect for public preferences, the upshot is that Iterated Backward Inference can be used to compute properly rationalizable strategies: if the procedure eliminates a strategy s_i , then s_i is not properly rationalizable.

6 An Iterated Elimination Procedure for Respect for Public Preferences

My iterated elimination procedure for deriving consequences of Respect for Public Preferences combines two principles: “local” dominance at an information set, and backward inference. We will see that these two principles are consequences of lexicographic rationality. I begin with dominance at an information set.

6.1 Sequential Admissibility

To compare two strategies s_i, s'_i at an information set I_i , I compare their payoffs against strategies s_{-i} that reach I_i . To ensure comparability, I require that s_i and s'_i are consistent with exactly the same strategies s_{-i} at I_i (in symbols, $\text{cons}(s_i, I_i) = \text{cons}(s'_i, I_i)$). In Section 9 we shall see that in games with perfect recall, $\text{cons}(s_i, I_i) = [I_i]_{-i}$ for any strategy s_i that reaches I_i . Hence in games with perfect recall, the condition that $\text{cons}(s_i, I_i) = \text{cons}(s'_i, I_i)$ can be replaced by the condition that s_i and s'_i are both consistent with I_i . This yields the following definition of dominance at an information set.

Definition 11 *Let T be a game tree with perfect recall with information set I_i belonging to player i .*

1. s_i weakly dominates s'_i at $I_i \iff$
 - (a) s_i and s'_i are each consistent with I_i (i.e., $s_i, s'_i \in [I_i]_i$), and
 - (b) there is a strategy s_{-i} consistent with I_i (i.e., $s_{-i} \in [I_i]_{-i}$) such that $u_i(s_i, s_{-i}) > u_i(s'_i, s_{-i})$, and
 - (c) for all s_{-i} consistent with I_i , we have that $u_i(s_i, s_{-i}) \geq u_i(s'_i, s_{-i})$.
2. s_i strictly dominates s'_i at I_i given $\Sigma_{-i} \subseteq S_{-i}(T) \iff$
 - (a) s_i and s'_i are each consistent with I_i , and
 - (b) I_i is consistent with Σ_{-i} (i.e., $[I_i]_{-i} \cap \Sigma_{-i} \neq \emptyset$), and
 - (c) for all $s_{-i} \in [I_i]_{-i} \cap \Sigma_{-i}$, we have that $u_i(s_i, s_{-i}) > u_i(s'_i, s_{-i})$.

According to the first part of this definition, a strategy s_i weakly dominates s'_i at an information set I_i if s_i weakly dominates s'_i if we restrict the space of strategies of the other player $-i$ to those that reach information set I_i . The definition of strict dominance at an information set is similar to that of weak dominance, but has an extra parameter Σ_{-i} to represent the result of previous rounds of elimination. To illustrate, in the Centipede game (Figure 1) for player 1 the strategy lt weakly dominates ll at information set I^1 (and hence in the entire game), and strictly dominates ll at I^3 given $\{l_2, t_2\}$ (that is, no matter what player 2 chooses). The strategy tt strictly dominates both lt and ll at I^1

given $\{t_2\}$. For player 2, the strategy t_2 strictly dominates l_2 at I^2 given $\{lt\}$. In Kohlberg's game (Figure 2), for player II the strategy UL strictly dominates DL at information set I^1 given $S_I(T)$ (that is, no matter what player I does). At I^2 , the strategy yT strictly dominates nB given $S_{II}(T)$, and yT strictly dominates nT given $\{DR\}$. Finally, at I^4 the strategy nB strictly dominates nT given $\{DR\}$.

A simple but important fact is that strict dominance is robust in the sense that if a strategy s_i strictly dominates another strategy s'_i , given some restriction Σ_{-i} , dominance still obtains under a stricter restriction $\Sigma'_{-i} \subset \Sigma_{-i}$, provided of course that Σ'_{-i} is consistent, that is, not empty. The next lemma records this fact for future reference.

Lemma 12 *Let T be a game tree with perfect recall. Suppose that a strategy s_i strictly dominates s'_i given $\Sigma_{-i} \subseteq S_{-i}(T)$ at information I_i . Let $\Sigma'_{-i} \subseteq \Sigma_{-i}$ be such that $[I_i]_{-i} \cap \Sigma'_{-i} \neq \emptyset$. Then s_i strictly dominates s'_i at I_i given Σ'_{-i} .*

The proof is immediate from Definition 11.

6.2 Entailment Inference and Backward Inference

In addition to local dominance at an information set, the second component of the iterated elimination procedure draws inferences from the results of elimination at one information set I to eliminate strategies at another information set I' . Consider two information sets I, I' such that all strategies s_i for player i consistent with I are also consistent with I' . Let us say that in this case I **entails I' for player i** ; so I entails I' for player i iff $[I]_i \subseteq [I']_i$. The general principle is this: if a strategy s_i consistent with I is considered unlikely given the information in the set I , and I entails I' for player i , then s_i is considered unlikely given the information in the set I' . Intuitively, if s_i is considered unlikely given I , then there is a possibility s'_i consistent with I that is considered more likely than s_i . Since I entails I' , the possibility s'_i is consistent with I' and hence s_i is not among the most likely possibilities given I' . To give the principle a precise formulation, I interpret “ s_i is considered unlikely at information set I ” to mean that the revision on I rules out the strategy s_i . Then in symbols, the **entailment inference principle** is that

$$\text{if } [I]_i \subseteq [I']_i, \text{ and } s_i \in [I]_i, \text{ but } s_i \notin \rho_{-i} * [I]_i, \text{ then } s_i \notin \rho_{-i} * [I']_i. \quad (1)$$

Or contrapositively: If I entails I' for player i , then $(\rho_{-i} * [I']_i) \cap [I]_i \subseteq \rho_{-i} * [I]_i$. For example, in the Centipede game, if $lt \notin \rho_2 * [I^3]_1$, then the entailment inference principle implies that $lt \notin \rho_2 * [I^2]_1$. In Kohlberg's game, if $nT \notin \rho_I * [I^4]_{II}$, then the principle implies that $nT \notin \rho_I * [I^2]_{II}$. On the other hand, in Kohlberg's game it is not the case that $[I^1]_2 \subseteq [I^2]_2$ (since for example UL is consistent with I^1 but not with I^2), so even if DL is considered unlikely at I^1 , the principle does not allow us to infer that DL is considered unlikely at I^2 . For example, we may have that $\rho_1 * [I^1]_2 = \{UL, UR\}$ and

$\rho_1 * [I^2]_2 = \{DL, DR\}$. Thus Principle 1 does not in general license “forward induction” arguments.¹

Lexicographic revision satisfies Principle 1. The next Lemma shows that this is due to a very general property of lexicographic probability systems which holds independent of a particular game structure.

Lemma 13 *Let ρ be an LPS with full support, and suppose that $P \subseteq P'$. Then $(\rho * P') \cap P \subseteq \rho * P$.*

Proof. Let $x' \in (\rho * P') \cap P$. Then x' is minimal in P' , that is, $x^* \leq_\rho x'$ for all $x^* \in P'$. Since $P \subseteq P'$, it follows that x' is minimal in P . Hence $x' \in \rho * P$. ■

If we set $[I]_i = P$, and $[I']_i = P'$, it is apparent that Principle 1 is an instance of (the contrapositive of) Lemma 13. Standard backward inference can be seen as an instance of Principle 1. Intuitively, in backward inference a player “looks ahead” and “reasons back”. Suppose that we have an information set I_i belonging to player i that follows some other information set I , which may belong to player $-i$ (in symbols, $I_i \geq I$). In Clause 3 of Lemma 25 in Section 9 below I show that in this case I_i entails I for player i (given perfect recall). So by Principle 1, we have that if a strategy s_i consistent with I is considered unlikely given the information in the set I_i , then it is considered unlikely given the information in the set I . In symbols, we have that for *games with perfect recall*:

$$\text{if } I_i > I, \text{ and } s_i \in [I]_i, \text{ but } s_i \notin \rho_{-i} * [I]_i, \text{ then } s_i \notin \rho_{-i} * [I]_i. \quad (2)$$

I refer to the instance 2 of Principle 1 as the **backward inference principle**. Given that $[I]_i \subseteq [I']_i$, Principle 2 is an instance of Principle 1. Note that the backward inference principle generalizes standard backward induction for games of perfect recall but imperfect information. This is especially clear when we have two information sets I_i, I_{-i} such that $I_i > I_{-i}$ and $I_{-i} > I_i$ (which is impossible in games of perfect information, but possible in games with perfect recall). In this case neither I_i nor I_{-i} is “lower than” the other, so traditional backward induction does not apply, but Principle 2 does.

Backward inference is the typical application of Principle 1. A different case in which the principle applies occurs when we have two information sets I, I' such that $I < I'$, and *all* strategies consistent with information set I are consistent with I' —that is, $[I] \subseteq [I']$, and hence $[I]_i \subseteq [I']_i$. This special case arises in a game with just two information sets, such as results from transcribing a game matrix directly into a game tree. For example, Figure 3 shows a game tree in which players’ options are just their strategies in the Centipede game. If for example $l_2 \notin \rho_1 * [I^1]_2$, then $l_2 \notin \rho_1 * [I^2]_2$; in general, we have that $\rho_1 * [I^2]_2 \subseteq \rho_1 * [I^1]_2$.

¹The following principle comes close to forward induction arguments: if $s_i \in [I]_i - \rho_{-i} * [I]$, and $\rho_{-i} * [I]_i \cap [I']_i \neq \emptyset$, then $s_i \notin \rho_{-i} * [I']_i$. This implies that if $\rho_1 * [I^1]_2 = \{UL, UR, DR\}$, then $\rho_1 * [I^2]_2 = \{DR\}$ as required by forward induction. Lexicographic revisions satisfy this principle; note however that it does not depend on I or I' being ordered in any particular way.

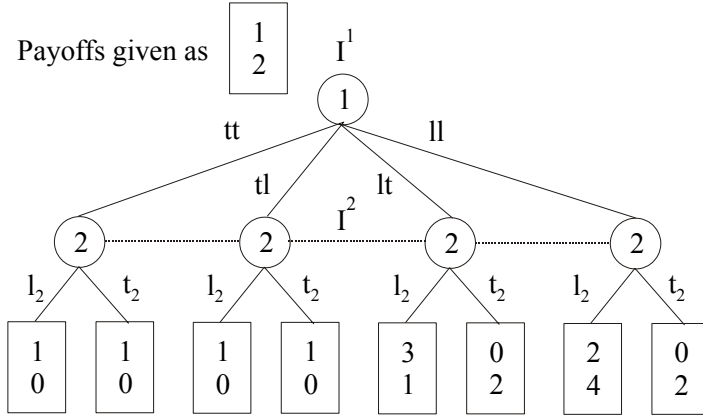


Figure 3: A game tree corresponding to the strategic form of the Centipede game

6.3 Iterated Backward Inference: Definition and Examples

Now I define the iterated elimination procedure, which I refer to as Iterated Backward Inference, or IBI for short.

Definition 14 (Iterated Backward Inference) *Let T be a game tree. Then define for all players i , for all information sets I_i , strategies $s_i \in S_i(T)$, $s_{-i} \in S_{-i}(T)$:*

1. $s_i \in \Gamma_i^0(I) \iff$
 - (a) s_i is consistent with I (i.e., $s_i \in [I]_i$), and
 - (b) for all information sets I_i belonging to player i such that $[I]_i \subseteq [I]_i$, we have that s_i is not weakly dominated at I_i given $S_{-i}(T)$.
2. $s_i \in \Gamma_i^{n+1}(I) \iff$
 - (a) $s_i \in \Gamma_i^n(I)$, and
 - (b) for all information sets I_i belonging to player i such that $[I]_i \subseteq [I]_i$, we have that s_i is not strictly dominated at I_i given $\Gamma_{-i}^n(I_i)$.

We may think of round n of Iterated Backward Inference as assigning a set of uneliminated strategies $\Gamma_i^n(I)$ and $\Gamma_{-i}^n(I)$ for each player, to each information set I . To begin with, if strategy s_i is not eliminated at round n at I , then s_i must be consistent with I . Furthermore, in round 0, s_i is eliminated at an information set I if there is another information set I_i that entails I for i , such that s_i is weakly dominated at I_i . In a later round $n + 1$, a strategy s_i is

information set/ surviving strategies	I^1	I^2	I^3
Γ_1^0	tt, tl, lt	lt	lt
Γ_2^0	l_2, t_2	l_2, t_2	l_2
Γ_1^1	tt, tl, lt	lt	lt
Γ_2^1	t_2	t_2	l_2
Γ_1^2	tt, tl	lt	lt
Γ_2^2	t_2	t_2	l_2

Table 7: Iterated Backward Inference in the Centipede Game (see Figure 1)

eliminated at I if there is an information set I_i that entails I for i , such that s_i is strictly dominated at I_i given the result of previous rounds of elimination.

Iterated Backward Inference assigns a set of strategies to an information set I rather than computing predictions for the game tree T as a whole. However, we may take the strategies that survive Iterated Backward Inference at *every* information set at stage n to be those that survive IBI for the game as a whole at stage n . Accordingly, define $\Gamma_i^n(T) =_{df} \{s_i : \forall I_i. \text{ if } s_i \in [I_i]_i, \text{ then } s_i \in \Gamma_i^n(I_i)\}$. The set of strategy profiles consistent with IBI at stage n is then given by $\Gamma^n(T) = \Gamma_1^n \times \Gamma_2^n$. It is immediate from the definition of the Γ procedure that at stage $n + 1$, the strategy set surviving for each player i at each information set I becomes smaller or stays the same. In a finite game, in which there are only finitely many strategies and information sets, this entails that eventually no more strategies will be eliminated at any information set. Let $\max(T)$ be the least stage m at which no additional strategy is eliminated at any information set. I write $\Gamma_i^\infty(I)$ for $\Gamma_i^{\max}(I)$, and similarly $\Gamma^\infty(T)$ for $\Gamma^{\max}(T)$. Similarly I write $\Gamma_i^\infty(I)$ for $\Gamma_i^{\max}(I)$, and similarly $\Gamma^\infty(T)$ for $\Gamma^{\max}(T)$. The upshot is that we may treat $\Gamma^\infty(I)$ as the prediction of Iterated Backward Inference for play reaching the information set I , and treat $\Gamma^\infty(T)$ as the prediction of Iterated Backward Inference for the game tree T .

To illustrate the procedure, I show its computations on the games of Figures 1, 3 and 2. The columns correspond to information sets in the game and the rows show the set of strategies $\Gamma_i^n, \Gamma_{-i}^n$ surviving at each round of elimination, for each player. Table 7 shows that the final result of the computation in the Centipede game is for Player 1 to take at all his information sets, and for player 2 to take at his. Formally, we have that $\Gamma_1^\infty(T) = \{tt, tl\}$ and $\Gamma_2^\infty(T) = \{t_2\}$. This result agrees with standard backward induction; in Section 8 I establish that Iterated Backward Inference agrees with backward induction in all perfect information games with a unique subgame-perfect equilibrium, of which the Centipede game is an example.

Table 8 shows that Iterated Backward Inference eliminates just one strategy in the matrix tree (normal form) of the Centipede game, namely ll which is weakly dominated by lt . After ll is dominated, there is no strict dominance, and the procedure terminates.

information set/ surviving strategies	I^1	I^2
Γ_1^0	tt, tl, lt	tt, tl, lt
Γ_2^0	l_2, t_2	l_2, t_2

Table 8: Iterated Backward Inference in a matrix game tree for the Centipede game (see Figure 3)

information set/ surviving strategies	I^1	I^2	I^3	I^4
Γ_1^0	yT, yB, nT	yT, yB, nT	nT, nB	nT, nB
Γ_2^0	UL, UR, DR	DL, DR	DL, DR	DL, DR
Γ_1^1	yT, yB, nT	yT, yB, nT	nT, nB	nT, nB
Γ_2^1	UL, UR	DL, DR	DL, DR	DL, DR

Table 9: Iterated Backward Inference in Kohlberg's game (Figure 2)

This example illustrates two main points. First, in games with two information sets, which are essentially just another representation of the strategic form of a game, Iterated Backward Inference coincides with the Dekel-Fudenberg procedure: First eliminate all weakly dominated strategies, then iteratively eliminate strictly dominated strategies. In light of Theorem 15 below, it follows that the Dekel-Fudenberg procedure is valid for deriving consequences of Respect for Public Preferences. The second point is that Iterated Backward Inference can yield different results for game trees that have the same normal form, as the two game trees for the Centipede Game do. In each case, the output of Iterated Backward Inference is valid in the sense that Respect for Public Preferences entails that eliminated strategies will not be played. In our examples, applying IBI in the game tree of Figure 3 yields the result that Player 1 does not choose ll , and applying IBI in the game tree of Figure 1 yields the result that Player 1 chooses neither ll nor lt . This illustrates how in some game trees IBI provides more information than in others, in that the procedure finds more of the consequences of Respect for Public Preferences. Thus although my main *epistemic principle*, Respect for Public Preferences, pertains to the strategic form of a game, and hence is independent of any particular extensive form, the *computational procedure* that I investigate in this paper does depend on a particular choice of game tree. For a given game G in strategic form, we would like to have a *canonical* game tree $T(G)$ such that a strategy s_i is consistent with Respect for Public Preferences just in case s_i survives IBI in $T(G)$. The most important formal question left open in this paper is the existence and construction of canonical game trees for IBI.

Table 9 shows the computation of Iterated Backward Inference in Kohlberg's game. This example illustrates two points. First, even though Kohlberg's game is not a game of perfect information, Iterated Backward Inference yields a unique

outcome prediction: that player II will choose U immediately, resulting in pay-offs $(0, 2)$. (For $\Gamma_{II}^\infty(T) = \{UL, UR\}$.) Second, the game shows how Iterated Backward Inference incorporates backward reasoning but not forward reasoning. At information set I^2 , both yT and yB strictly dominate nB , so nB is eliminated at I^2 at round 0, and hence by backward inference, nB is also eliminated at I^1 at round 0. After nB is eliminated, UL and UR strictly dominate DR at information set I^1 in round 1, leaving U as the only choice for player II . By contrast, there are two forward induction arguments that IBI does not incorporate. First, since DL is eliminated at I^1 , one might take forward induction to imply that DL should be eliminated at I^2 as well. Second, since nB is eliminated at I^2 , one might take forward induction to imply that nB should be eliminated at I^3 as well.

7 Soundness of Iterated Backward Inference and Existence of Solution

This section contains the main theorems of this paper. First, I show my key result: For finite games with perfect recall, if Respect for Public Preferences obtains and lexicographic rationality (with full support) is common belief, then it is common belief that each player believes that play follows Iterated Backward Inference, at each information set I . Second, I establish existence: in finite games there are some predictions consistent with Iterated Backward Inference.

7.1 Soundness

Theorem 15 *Let T be a finite game tree with perfect recall and assume that Lexicographic Rationality and Full Support are common belief (see Definition 6). Then Respect for Public Preferences (Axiom 9) implies that for all n, i, I :*

1. *if a strategy s_i is strictly dominated at an information set I given $\Gamma_{-i}^n(I)$, then $CB(\exists s'_i \in [I]_i. s'_i \succ_i s_i)$, and*
2. $CB(\rho_{-i} * [I]_i \subseteq \Gamma_i^n(I))$.

Recall that $choice_i$ denotes the strategy choice of player i . The strategy profile chosen is then given by $choice =_{df} choice_1 \times choice_2$. It is a simple corollary from Theorem 15 that Iterated Backward Inference makes correct predictions about the choices of the players, given Respect for Public Preferences and the standard epistemic assumptions.

Corollary 16 *Let T be a finite game tree with perfect recall and assume that Lexicographic Rationality and Full Support are common belief. Then Respect for Public Preferences (Axiom 9), Preference Maximization and Preference Introduction (see Definition 6) imply that $choice \in \Gamma^n(T)$ for all n .*

7.2 Existence of Solution

At a given information set I , the set of possible strategy combinations consistent with Iterated Backward Inference is given by $\Gamma^n(I) =_{df} \Gamma_1^n(I) \times \Gamma_2^n(I)$. The next lemma says that under perfect recall, the set of strategy profiles consistent with IBI decreases as we move down the game tree from one information set I_i to another I'_i .

Lemma 17 *Let T be a game tree with perfect recall. Then for all players i , information sets I_i, I'_i belonging to i , if $I'_i \geq I_i$, then $\Gamma^n(I_i) \cap [I'_i] \subseteq \Gamma^n(I'_i)$.*

The inclusion is not an equality because it may be the case that $\Gamma^n(I_i) \cap [I'_i] = \emptyset$. For example, in the Centipede game we have that $\Gamma^1(I^1) = \{(tt, t_2), (tl, t_2), (lt, t_2)\}$, and $\Gamma^1(I^3) = \{(lt, t_2)\}$; here $[I^3] = \{(ll, l_2), (lt, l_2)\}$ so $\Gamma^1(I^1) \cap [I^3] = \emptyset$. Given Lemma 17, we can establish that in finite games with perfect recall, Iterated Backward Inference returns a nonempty result.

Proposition 18 *Let T be a finite game tree with perfect recall. Then $\Gamma^n(T) \neq \emptyset$ for all n .*

It is easy to see that in infinite games existence may fail. The simplest example is a 1-player game in which the player may choose any natural number k , and the utility function is just $u_1(k) = k$. Since each option is strictly dominated in this game, IBI eliminates all options. For finite games, we have the following characterization of strategies that survive all rounds of elimination.

Lemma 19 *Let T be a finite game tree with perfect recall. A strategy s_i is in $\Gamma_i^\infty(I) \iff$ for all information sets I_i such that s_i is consistent with I_i and I_i entails I for I :*

1. s_i is admissible at I_i given $S_{-i}(I_i)$, and
2. s_i is not strictly dominated at I_i given $\Gamma_{-i}^\infty(I_i)$.

8 Backward Induction

Backward inference is more general than traditional backward induction in that it applies to games with imperfect information (cf. Section 6.2). In this section I show that in finite perfect information games or repeated stage games in which backward induction yields a unique solution, Iterated Backward Inference agrees with the backward induction solution. In such games Iterated Backward Inference includes backward induction as a special case.

8.1 Preliminaries

A game tree T has **perfect information** just in case each information set in T is a singleton. For game trees with perfect information I write x instead of

$\{x\}$ to denote an information set. In a game of perfect information, each node x is the root of a subgame of perfect information, which I denote by T_x . I write $u_i(s_i, s_{-i}, x)$ for the payoff to player i that results if the game begins at node x and follows strategies s_i and s_{-i} . A strategy pair (s_i, s_{-i}) is a **subgame perfect equilibrium** in T iff for each player i , we have $u_i(s_i, s_{-i}, x) \geq u_i(s'_i, s_{-i}, x)$ for all strategies s'_i and all nodes x in T . If a game tree has a *unique* subgame perfect equilibrium, I write $b_i(x)$ for the payoff that player i receives in the subgame T_x in the unique subgame perfect equilibrium. Note that a game tree with perfect information has a unique subgame perfect equilibrium just in case each subgame T_x has a unique subgame perfect equilibrium.

For example, in the Centipede game of Figure 1, the unique subgame perfect equilibrium is (tt, t_2) , and $b_1(\emptyset) = 1, b_1(\text{leave}) = 0, b_1(\text{leave} * \text{leave}) = 3$, while $b_2(\emptyset) = 0, b_2(\text{leave}) = 2, b_2(\text{leave} * \text{leave}) = 1$.

Suppose that player i moves first in a game tree T with root r , and that T has a unique subgame perfect equilibrium (s_i, s_{-i}) . I say that move a for player i is a BI-maximizer iff $b_i(r * a) \geq b_i(r * a')$ for all other moves $r * a'$ in T . It is easy to see that $s_i(r)$ must be a BI-maximizer.

Lemma 20 *Let T be a game tree of perfect information with a unique subgame perfect equilibrium (s_i, s_{-i}) . Then a strategy profile (s'_i, s'_{-i}) is equal to $(s_i, s_{-i}) \iff$*

1. $s'_i(r)$ is a BI-maximizer, and
2. (s'_i, s'_{-i}) is an SPE in each subtree T_{r*a} .

8.2 Iterated Backward Inference and Backward Induction

The main difference between Iterated Backward Inference and subgame perfection is that IBI does not consider the behaviour of a strategy s_i at an information set I_i that is unreachable with s_i . For example, in the Centipede game, the strategy profile (tl, t_2) survives IBI (i.e., $(tl, t_2) \in \Gamma^\infty(T)$), but tl, t_2 is not a subgame perfect equilibrium because tl chooses l at the node $\text{leave} * \text{leave}$, which is not a BI-maximizer. However, I show that if s_i survives IBI, then at all *reachable* information sets, s_i agrees with subgame perfection (backward induction). I say that a profile (s_i, s_{-i}) is **extendible to an SPE** (subgame perfect equilibrium) in a tree T iff there is a strategy profile (s'_i, s'_{-i}) such that

1. (s'_i, s'_{-i}) is an SPE in T , and
2. for all x in T : if $\text{player}(x) = i$ and $s_i \in [x]_i$, then $s_i(x) = s'_i(x)$, and likewise for player $-i$.

For example, in the Centipede game (tt, t_2) extends (tl, t_2) . I observe that if (s'_i, s'_{-i}) extends (s_i, s_{-i}) , then $\text{play}(s_i, s_{-i}) = \text{play}(s'_i, s'_{-i})$. In terms of extendibility, the fact that a strategy s_i surviving IBI agrees with backward induction at any information set I_i consistent with s_i amounts to the requirement that

s_i can be extended to a subgame perfect strategy s'_i by setting $s_i(I_i) = s'_i(I_i)$ whenever I_i is inconsistent with s_i . The next proposition verifies that this is the case for strategies surviving Iterated Backward Inference.

Proposition 21 *Let T be a finite game tree with perfect information and a unique subgame perfect equilibrium. Then for each node x , for each strategy profile $(s_i, s_{-i}) \in \Gamma^\infty(x)$, we have that*

1. (s_i, s_{-i}) is extendible to an SPE in T_x , and
2. $u_i(s_i, s_{-i}, x) = b_i(x)$.

Proposition 21 implies that the play consistent with IBI is exactly the backward induction play. For if r is the root of the game tree T , then it is easy to see that $\Gamma^\infty(r) = \Gamma^\infty(T)$ —a strategy survives IBI just in case it survives each round of elimination at the root. So by Proposition 21, every strategy pair $(s_i, s_{-i}) \in \Gamma^\infty(T)$ is extendible to the unique SPE in $T_r = T$. So the play sequence of every strategy pair surviving IBI is the backward induction path.

In addition to game trees with perfect information, subgame perfection yields a unique outcome in a finite repetition of a stage game with a unique Nash equilibrium. A well-known example is the Prisoner's Dilemma. We can show that if IBI yields a unique outcome prediction for a game G , its prediction for the repeated game G^k is that each repetition will have the same outcome as predicted for the stage game. Thus in a finitely iterated Prisoner's Dilemma, IBI predicts defection at each stage, as backward induction does. The demonstration of this observation is similar to the proof for Proposition 21; I sketch the definitions and omit the proofs.

Say that a game G is **solvable by Iterated Backward Inference** just in case for all strategy profiles (s_i, s_{-i}) and $(s'_i, s'_{-i}) \in \Gamma^\infty(G)$, and each player i , we have $u_i(s_i, s_i) = u_i(s'_i, s'_{-i})$. Let G^k be the k -repetition of G . I say that an information set I in G^k is at stage m if $m - 1$ repetitions of the game have been played before I . Then we have the following proposition.

Proposition 22 *Let $s_i \in \Gamma^\infty(G^k)$. For all stages $m \leq k$, for all information sets I , if I is at stage m , then $s_i(G_m) \in \Gamma^\infty(G)$.*

For further discussion of the relationship between backward induction, respect for preferences and proper rationalizability see [16] and [1].

9 Perfect Recall: Basic Lemmas and the Combination Principle

This section establishes some lemmas about game trees with perfect recall that relate the sequential information structure of a game tree to what players know about their opponents' strategies given the information sets in the game. The main result in this section is to establish that the following fact—which I call the

Combination Principle—holds under perfect recall: Given two strategies s_i, s'_i for player i that both reach an information set I_i , there is a third strategy s_i^* that agrees with s_i on all play sequences that reach I_i , and that agrees with s'_i on all play sequences that do not reach I_i .

For a given game tree T , I follow [13, p.203] and define the **epistemic history** $E_i(x)$ for player i at node x as follows.²

1. $E_i(\emptyset) = \emptyset$.
2. $E_i(x * a) = E_i(x) * I(x) * a$ if $\text{player}(x) = i$, and
3. $E_i(x * a) = E_i(x)$ if $\text{player}(x) \neq i$.

Thus the epistemic history of player i at node x contains the information that player i learned before reaching node x . For example, in the Kohlberg game $E_2(D * n * T) = E_2(D) = \langle I^1, D \rangle$.

Definition 23 (Osborne and Rubinstein) *A game tree T has **perfect recall** if for each player i , for all information sets I_i , for all nodes $x, x' \in I_i$, we have $E_i(x) = E_i(x')$.*

I now establish a number of lemmas that relate the sequential structure of information sets to differences in player's knowledge at different information sets; they can be skipped without loss of continuity. Proofs are in Section 13.

Lemma 24 *Let T be a game tree with perfect recall. Then*

1. for all nodes x, x' if $x < x'$ then $I(x) \neq I(x')$.
2. if $s_i \in [I_i]_i$ and $(s'_i, s_{-i}) \in I_i$, then $I_i \cap \text{Play}(s_i, s_{-i}) = I_i \cap \text{Play}(s'_i, s_{-i})$.
3. if $(s_i, s_{-i}) \in [I_i]$ and $(s'_i, s'_{-i}) \in [I_i]$, then $(s_i, s'_{-i}) \in [I_i]$.

The next lemma mainly concerns the \geq ordering among information sets.

Lemma 25 *Let T be a game tree with perfect recall. Let I_i be an information set belonging to player i . Then*

1. for all $s_i \in [I_i]_i$ we have that $\text{cons}(s_i, I_i) = [I_i]_{-i}$.
2. for all information sets I'_i , if $I_i \geq I'_i$, then $[I_i] \subseteq [I'_i]$.
3. for all information sets I , if $I_i \geq I$ then $[I_i]_i \subseteq [I]_i$.
4. for all information sets I_i, I'_i of player i it is the case that: $I_i = I'_i \iff (I_i \geq I'_i \text{ and } I'_i \geq I_i)$.

²Osborne and Rubinstein use a slightly different definition which, however, yields an equivalent notion of perfect recall.

Lemmas 24 and 25 allow us to establish the following Combination Principle. Let s_i, s'_i be two strategies for player i . I say that s_i^* agrees with s'_i on I_i if for all strategies s_{-i} , if $(s_i, s_{-i}) \in [I_i]$, then $play(s_i^*, s_{-i}) = play(s'_i, s_{-i})$. Similarly, s_i^* agrees with s_i outside of I_i if for all strategies s_{-i} , if $(s_i, s_{-i}) \notin [I_i]$, then $play(s_i^*, s_{-i}) = play(s_i, s_{-i})$. The Combination Principle says that if s_i, s'_i are each consistent with an information set I_i , then s_i and s'_i can be combined to yield a third strategy s_i^* that agrees with s_i on I_i and agrees with s'_i outside of I_i . Intuitively, the following instructions define the combined strategy s_i^* : First, follow s_i until either information set I_i is reached or play arrives at an information set I'_i from which I_i is unreachable. In the latter case, follow s'_i . In the former case, follow s_i . These instructions require an agent to remember whether I_i has been reached or not—hence the importance of perfect recall. The next proposition is a formal statement of the Combination Principle.

Proposition 26 (Combination Principle) *Let T be a game tree with perfect recall. Let s_i, s'_i be two strategies consistent with an information set I_i (i.e., $s_i, s'_i \in [I_i]_i$). Then there is a strategy s_i^* such that*

1. s_i^* agrees with s'_i on I_i (i.e., for all strategies s_{-i} if $(s'_i, s_{-i}) \in [I_i]$, then $play(s_i^*, s_{-i}) = play(s'_i, s_{-i})$), and
2. s_i^* agrees with s_i outside of I_i (i.e., for all strategies s_{-i} , if $(s_i, s_{-i}) \notin [I_i]$, then $play(s_i^*, s_{-i}) = play(s_i, s_{-i})$).

10 Lexicographic Rationality and Sequential Admissibility

A common principle of rational choice in extensive form game holds that a rational strategy should be sequentially rational, that is, rational at each information set. As Blume, Brandenburger and Dekel point out, lexicographic rationality satisfies this principle to a considerable extent: if ρ is an LPS with full support that maximizes lexicographic expected utility, then it follows from Proposition 4 that ρ must maximize lexicographic expected utility on each “information cell” [3, p.61/62]. From the point of view of individual decision theory, this is a more or less immediate consequence of the sure thing principle or the Independence Axiom of decision theory. In games that satisfy the Combination Principle from the previous section, the Independence Axiom implies that strategies that maximize lexicographic expected utility must do so at each information set; that is, they must be sequentially rational. For suppose that some strategy s_i fails to maximize lexicographic expected utility given the information in some information set I_i ; let s'_i be preferred to s_i given I_i . Then by the Combination Principle, there is a third strategy s_i^* that agrees with s'_i on I_i and with s_i outside of I_i . Thus s_i^* yields higher expected lexicographic utility than s_i conditional on I_i and yields the same outcomes as s_i outside of I_i . So it follows from the sure thing principle (Proposition 4) that s_i^* is preferred to s_i . To illustrate,

in the Centipede game, the strategy lt for player 1 is strictly preferred to ll at information set I^3 . On all play sequences that don't reach I^3 (i.e., if player 2 doesn't choose "leave"), the strategies lt and ll yield the same payoff. Hence if player 1 maximizes lexicographic expected utility (with full support), it follows that player 1 prefers lt to ll . This is not the same as *assuming* that player 1 is "substantively rational", meaning that she would maximize expected payoffs at every information set that the play might arrive at. Instead I *derive* this fact from the assumption of lexicographic rationality, which pertains to the strategic form of the game, not its normal form.

The general connection between dominance and lexicographic rationality is this: if in a game of perfect recall, a strategy s_i is weakly dominated at an information set I_i , or strictly dominated given the revision $\rho_i * [I_i]_{-i}$, then s_i does not maximize lexicographic expected utility. For by the Combination Principle, there is a strategy s'_i that is preferred "locally" at information set I_i , and that behaves the same as s_i outside of I_i . By the sure thing principle, s'_i is then preferred to s_i . The next proposition formalizes this observation.

Proposition 27 *Let T be a finite game tree with perfect recall. Let ρ_i be an LPS for player i (i.e., ρ_i defines the preferences of player i over strategies) with full support. Consider any information set I_i . Suppose that (1) s_i is weakly dominated at I_i , or (2) s_i is strictly dominated at I_i given Σ_{-i} and $\rho_i * [I_i]_{-i} \subseteq \Sigma_{-i}$. Then there is a strategy s'_i such that*

1. $s'_i \in [I_i]_i$ and
2. $s'_i \sim^{|[I_i]_{-i}|} s_i$, and
3. $s'_i \succ^{|[I_i]_{-i}|} s_i$, and
4. $s'_i \succ_i s_i$.

11 Conclusion

An important approach to developing and understanding solution concepts for game theory is to examine the epistemic assumptions that underlie predictions about the outcome of a game. In this paper I considered the consequences of Respect for Public Preferences: if it is common belief that an agent A prefers option a to option b , then it is common belief that in a binary choice between a and b , the agent A chooses a . Following previous work by Blume *et al.* [4] and Asheim [1], I propose to capture Respect for Public Preferences by requiring that each agent $B \neq A$ assigns probability 1 to A 's choosing a , conditional on A choosing a or b , whenever A 's preference for a over b is common belief. We can employ lexicographic probability systems to ensure that this conditional probability is well-defined.

Iterated Backward Inference (IBI) is a procedure for computing the consequences of common belief in Revealed Preference in a given game G . Iterated

Backward Inference eliminates strategies in a game tree T . My main result is that the procedure is valid given common belief in Revealed Preference, in the following sense: if T is an extensive form of the strategic game G , and IBI eliminates a strategy s , then s is not chosen in the game G . Iterated Backward Inference generalizes two well-known algorithms for solving games: the Dekel-Fudenberg procedure (first eliminate weakly dominated strategies, then iteratively strictly dominates ones), and standard backward induction for game trees with perfect information; IBI yields predictions that are at least as strong as those given by these two algorithms.

It follows from Asheim’s characterization of proper rationalizability [1] that properly rationalizable strategies are consistent with Respect for Public Preferences. Hence IBI can be used to find strategies that are not properly rationalizable.

IBI is valid for computing consequences of Respect for Public Preferences because of two key facts. First, given perfect recall, lexicographic rationality enforces sequential admissibility (admissibility at reachable information sets). Second, lexicographic rationality satisfies the entailment inference principle: Consider two information sets I, I' such that all strategies consistent with I are consistent with I' . Then if a strategy s is considered unlikely at I , the strategy s is also considered unlikely at I' . In the case in which the information set I is a successor of I' , the entailment inference principle yields a backward inference principle.

I mention two open questions for future research. In different game trees with the same strategic form G , Iterated Backward Inference may give stronger results. We would like to apply IBI to a canonical game trees $T(G)$ in which the procedures gives complete results, eliminating all and only those strategies inconsistent with Respect for Public Preferences. Whether canonical game trees exist for an arbitrary game and how to construct them is perhaps the most important open formal question for understanding the computational aspects of Respect for Public Preferences, and perhaps proper rationalizability as well. The example of the game trees for the Centipede game (Figures 1 and 3) suggest that canonical game trees are those that in some sense have “as many information sets as possible” consistent with the given strategic form.

Respect for Public Preferences does not validate typical forward induction arguments (for example, it does not entail the forward induction solution in the well-known Burning Dollar game). We would like to know further epistemic principles that underlie forward induction arguments. Is there a “best rationalization” principle based on Respect for Public Preferences that validates forward induction?

12 Acknowledgements

I am indebted to Geir Asheim, Mamoru Kaneko, Cristina Bicchieri and Phil Curry for helpful discussions. Anonymous referees for the TARK IX conference suggested improvements to the paper. This research was supported by a grant

from the Social Sciences and Humanities Research Council of Canada.

13 Proofs

I introduce some additional notation to facilitate the proofs. It is useful to write $I(x) = I_i$ to express the fact that finite history x belongs to the cell I_i of player i 's information partition. For a finite action sequence $x = (a_1, \dots, a_n, \dots)$, let $move(x, n) =_{df} a_n$. If a node x precedes a node x' (i.e., $x < x'$), then there is a unique move a at x such that $x * a \leq x'$; I denote this move by $move(x, x')$. For example, in the Centipede Game, $move(\emptyset, leave * leave * leave) = leave$.

Since a strategy s_i assigns the same move to all nodes in the same information set, a strategy can be viewed as a function of the information sets of player i as well as a function of the nodes belonging to player i ; I write both $s_i(x)$ and $s_i(I(x))$ depending on what is most concise in context. Similarly, under perfect recall, all nodes in an information set I share the same epistemic history; in that case I will sometimes write $E_i(I)$ for the shared epistemic history at information set I .

It is sometimes useful to consider a history not as a sequence of actions but as a set of finite action sequences; therefore I write $Play(x) =_{df} \{x' : x' \leq x\}$ and similarly $Play(h) =_{df} \{x : x < h\}$ for an infinite history h . To shorten notation, let $Play(s_1, s_2) =_{df} Play(play(s_1, s_2))$. For example, in the Centipede game, the set of nodes reached during the play of $\langle l, t \rangle$ and t_2 is given by $Play(\langle l, t \rangle, t_2) = \{\emptyset, \emptyset * leave, \emptyset * leave * take\}$.

Lemma 10 *Common Belief in Axiom 7 implies Axiom 9.*

Proof. If Axiom 7 is common belief, then (a) $CB(B_i(s_{-i} \succ s'_{-i}) \rightarrow \rho_i * \{s_{-i}, s'_{-i}\}) = \{s_{-i}\}$. Since Lemma 8 is a theorem, it is common belief and hence we have (b) $CB([\rho_i * \{s_{-i}, s'_{-i}\} = \{s_{-i}\}] \rightarrow [s_i >_{-i} s'_{-i}])$. Given that common belief is closed under implication, (a) and (b) imply that (c) $CB([B_{-i}(s_i \succ_i s)] \rightarrow [s_i >_{-i} s'_{-i}])$. Since common belief is closed under implication, it follows that (d) if $CB(B_{-i}(s_i \succ_i s))$, then $CB(s_i >_{-i} s'_{-i})$. We also have that (e) if $CB(s_i \succ_i s'_i)$, then $CB(B_{-i}(s_i \succ_i s'_i))$ from the definition of common belief. Combining (d) and (e) yields Axiom 9. ■

Theorem 15 *Let T be a finite game tree with perfect recall and assume that Lexicographic Rationality and Full Support are common belief (see Axiom 6). Then Respect for Public Preferences (Axiom 9) implies that for all n, i, I :*

1. *if a strategy s_i is strictly dominated at an information set I given $\Gamma_{-i}^n(I)$, then $CB(\exists s'_i \in [I]_i. s'_i \succ_i s_i)$, and*
2. $CB(\rho_{-i} * [I]_i \subseteq \Gamma_i^n(I))$.

Proof. Base Step, $n = 0$. Part 1: Consider an information set I_i with $player(I_i) = i$.

First I consider $\rho_{-i} * [I_i]_i$. I show that (a) for each s_i , if $s_i \notin \Gamma_i^0(I_i)$, then $CB(s_i \notin \rho_{-i} * [I_i]_i)$. Suppose that $s_i \notin \Gamma_i^0(I_i)$. If $s_i \notin [I_i]_i$, then it is immediate that $s_i \notin \rho_{-i} * [I_i]_i$. Otherwise there is a strategy \hat{s}_i and an information set I'_i such that (1) $[I'_i]_i \subseteq [I_i]_i$ and (2) \hat{s}_i weakly dominates s_i at I'_i given $S_{-i}(T)$. By Proposition 27, there is s'_i such that s'_i is consistent with I'_i and $s'_i \succ_i s_i$. Since this is a logical consequence of the game structure, which is common belief, we have that $CB(s'_i \succ_i s_i)$ for some s'_i . By Axiom 9 it follows that (b) $CB(s'_i >_{-i} s_i)$. Since $s'_i \in [I'_i]_i$, it follows from (1), which is common belief as part of the game structure, that $s'_i \in [I_i]_i$. So $CB(\exists s'_i. s'_i >_{-i} s_i \text{ and } s'_i \in [I_i]_i)$, and thus $CB(s_i \notin \rho_{-i} * [I_i]_i)$. This establishes (a). Since (a) holds for all strategies s_i of player i , it follows that $\forall s_i \notin \Gamma_i^0(I_i)[CB(s_i \notin \rho_{-i} * [I_i]_i)]$. Hence³ $CB(\forall s_i \notin \Gamma_i^0(I_i)[s_i \notin \rho_{-i} * [I_i]_i])$, which is equivalent to $CB(\rho_{-i} * [I_i]_i \subseteq \Gamma_i^0(I_i))$ given the assumption that common belief is closed under implication.

Second, consider $\rho_i * [I_i]_{-i}$. Let $s_{-i} \in \rho_i * [I_i]_{-i}$. It is immediate that $s_{-i} \in [I_i]_{-i}$. I show that (*) for all information sets I_{-i} such that $[I_{-i}]_{-i} \subseteq [I_i]_{-i}$: if $s_{-i} \in [I_{-i}]_{-i}$, then $s_{-i} \in \Gamma^0[I_{-i}]_{-i}$. If $s_{-i} \in (\rho_i * [I_i]_{-i}) \cap [I_{-i}]_{-i}$, we can apply the entailment inference Principle 1 to conclude that $s_{-i} \in \rho_i * [I_{-i}]_{-i}$. In the first part of the argument I established that $\rho_i * [I_{-i}]_{-i} \subseteq \Gamma_{-i}^0(I_{-i})$. So $s_{-i} \in \Gamma_{-i}^0(I_{-i})$, which establishes (*). From (*) it follows that s_{-i} is not weakly dominated at any information set I_{-i} such that $[I_{-i}] \subseteq [I_i]_{-i}$. Hence $s_{-i} \in \Gamma_{-i}^0(I_i)$, and in general $\rho_i * [I_i]_{-i} \subseteq \Gamma_{-i}^0(I_i)$.

Part 2: Suppose that s_i is strictly dominated at some I'_i given $\Gamma_{-i}^0(I'_i)$ where I'_i entails I_i for i . By Proposition 27, it follows that there is $s'_i \in I'_i$ such that for all ρ_i , if $\rho_i * [I'_i]_{-i} \subseteq \Gamma_{-i}^0(I'_i)$, then $s'_i \succ_i s_i$. Since this conclusion depends only on the game structure, which is common belief, we have that $CB(\rho_i * [I'_i]_{-i} \subseteq \Gamma_{-i}^0(I'_i) \rightarrow s'_i \succ_i s_i)$. By Part 1 of the current theorem, $CB(\rho_i * [I'_i]_{-i} \subseteq \Gamma_{-i}^0(I'_i))$. Since CB is closed under implication, it follows that $CB(s'_i \succ_i s_i)$.

Inductive Step: Assume the hypothesis for n and consider $n + 1$.

Part 1: Consider information set I_i with $player(I_i) = i$.

Again we begin with $\rho_{-i} * [I_i]_i$ and show the contrapositive: if $s_i \notin \Gamma_i^{n+1}(I_i)$, then $CB(s_i \notin \rho_{-i} * [I_i]_i)$. Suppose that $s_i \notin \Gamma_i^{n+1}(I_i)$.

Case 1: $s_i \notin [I_i]_i$. As before, the conclusion is immediate.

Case 2: $s_i \in [I_i]_i - \Gamma_i^n(I_i)$. Then it follows from the inductive hypothesis that $CB(s_i \notin \rho_{-i} * [I_i]_i)$.

Case 3: $s_i \in [I_i]_i \cap \Gamma_i^n(I_i)$. Then there is some information set I'_i entailing I_i for i such that s_i is strictly dominated at I'_i given $\Gamma_{-i}^n(I'_i)$. By Part 2 of the inductive hypothesis, there is a strategy $s'_i \in [I'_i]_i$ such that $CB(\exists s'_i. s'_i \succ_i s_i)$. Moreover, since $[I'_i]_i \subseteq [I_i]_i$ we have that $s'_i \in [I_i]_i$. Hence by Axiom 9, we obtain that (d) $CB(\exists s'_i \in [I_i]_i. s'_i >_{-i} s_i)$. From Lemma 8, we have (e) $CB([\exists s'_i \in [I_i]_i. s'_i >_{-i} s_i] \rightarrow s_i \notin \rho_{-i} * [I_i]_i)$. Combining (e) and (d) yields that $CB(s_i \notin \rho_{-i} * [I_i]_i)$.

³Interchanging the order of CB and \forall is unproblematic because in a finite game, there are only finitely many strategies for each player and hence a statement of the form $\forall s_i. CB(p(s_i))$ is equivalent to a finite conjunction of the form $CB(p(s_i^1)) \wedge \dots \wedge CB(p(s_i^k))$, which is equivalent to $CB(p(s_i^1) \wedge \dots \wedge p(s_i^k))$ and hence to $CB(\forall s_i. p(s_i))$.

Using the same argument as in the base case, I conclude that $\rho_i * [I_i]_{-i} \subseteq \Gamma_{-i}^{n+1}(I_i)$.

Part 2: Same as in the base case. ■

Corollary 16 *Let T be a finite game tree with perfect recall and assume that Lexicographic Rationality and Full Support are common belief. Then Respect for Public Preferences (Axiom 9), Preference Maximization and Introspection (see Axiom 6) imply that $choice \in \Gamma^n(T)$ for all n .*

Proof. It suffices to show the contrapositive version: for each i , for all strategies s_i , we have that if $s_i \notin \Gamma_i^n(T)$, then $s_i \neq choice_i$.

Base Case, $n = 0$. Suppose that $s_i \notin \Gamma_i^0(I_i)$ for some information set I_i where $s_i \in [I_i]_i$. Then there is a strategy s_i^* and an information set I'_i entailing I_i for i such that s_i^* weakly dominates s_i . So by Proposition 27, there is a strategy s'_i such that $s'_i \succ_i s_i$. Preference Maximization (see Axiom 6) says that $\exists s'_i. s'_i \succ_i s_i \rightarrow s_i \neq choice_i$. Hence $s_i \neq choice_i$.

Inductive step, $n + 1$. Suppose that $s_i \notin \Gamma_i^{n+1}(T)$. If $s_i \notin \Gamma_i^n(T)$, the claim follows from the inductive hypothesis. Otherwise $s_i \in \Gamma_i^n(I_i) - \Gamma_i^{n+1}(I_i)$ for some I_i . Then s_i is strictly dominated at some I'_i given $\Gamma_{-i}^n(I'_i)$, where I'_i entails I_i for i . As in Part 2 of Theorem 15, there is s'_i such that $CB(\exists s'_i. s'_i \succ s_i)$. By Definition of common belief, we have $B_i(\exists s'_i. s'_i \succ s_i)$. So by Preference Introspection, $s'_i \succ s_i$ and so by Preference Maximization, $s_i \neq choice_i$. Since this is true for all i, n , we have that $\forall n. choice \subseteq \Gamma^n$. ■

Lemma 17 *Let T be a game tree with perfect recall. Then for all players i , information sets I_i, I'_i , if $I'_i \geq I_i$, then $\Gamma^n(I_i) \cap [I'_i] \subseteq \Gamma_{-i}^n(I'_i)$.*

Proof. If $I'_i \geq I_i$, then by Clause 3 of Lemma 25, we have that $[I'_i]_i \subseteq [I_i]_i$. I show the contrapositive of the consequent. Suppose that $s_i \notin \Gamma_i^n(I'_i)$. If $s_i \notin [I'_i]_i$, the claim follows immediately. So suppose that $s_i \in [I'_i]_i - \Gamma_i^n(I'_i)$.

Case 1: $n = 0$. Suppose that $s_i \notin \Gamma_i^0(I'_i)$. Then there is an information set I_i^* entailing I'_i for player i such that s_i is weakly dominated at I_i^* given $S_{-i}(T)$. Since $[I'_i]_i \subseteq [I_i^*]_i \subseteq [I_i]_i$, we have that I_i^* entails I_i for i , and so $s_i \notin \Gamma_i^0(I_i)$.

Case 2: $n > 0$. Suppose that $s_i \notin \Gamma_i^n(I'_i)$. Then there is an information set I_i^* such that s_i is strictly dominated at I_i^* given $\Gamma_{-i}^{n-1}(I_i^*)$ and I_i^* entails I'_i for i . As in the previous case, it follows that I_i^* entails I_i for i , and so $s_i \notin \Gamma_i^n(I_i)$. ■

Proposition 18 *Let T be a finite game tree with perfect recall. Then $\Gamma^n(T) \neq \emptyset$ for all n .*

Proof. Case 1: $n = 0$. In a finite game tree T , there is clearly a strategy s_i for player i that is not weakly dominated in T since weak dominance is transitive. If s_i is not weakly dominated in T , then s_i is not weakly dominated at any information set. For if s_i is weakly dominated at I_i , we may apply Case 1 of Proposition 27 to conclude that s_i is weakly dominated in T , contrary to supposition. Hence $s_i \in \Gamma_i^0(T)$.

Case 2: $n > 0$. For each information set I_i , there is unique history $E_i(I_i)$ shared by all nodes in I_i . Thus we may define $height_i(I_i) = |E_i(I_i)|$. Let $H = \max\{height_i(I_i) : I_i \text{ belongs to } i\}$. I note that if $I'_i > I_i$, then $height_i(I'_i) > height_i(I_i)$. For $n > 0$, say that s_i is n -acceptable at I_i if for all $I'_i \geq I_i$, the strategy s_i is not strictly dominated at I_i given $\Gamma_{-i}^{n-1}(I'_i)$. I argue by induction on $k = H - height(I_i)$ that if $H - height(I_i) = k$ there there is an n -acceptable strategy s_i at I_i .

Base Case, $k = 0$. Then $height(I_i) = H$, and so there is no I'_i such that $I'_i > I_i$. Suppose that s_i is strictly dominated at I_i given $\Gamma_{-i}^{n-1}(I_i)$. Since T is finite, and strict dominance is transitive, there clearly is a strategy s'_i that strictly dominates s_i given $\Gamma_{-i}^{n-1}(I_i)$ such that s'_i is not strictly dominated at I_i given $\Gamma_{-i}^{n-1}(I_i)$. So s'_i is n -acceptable at I_i .

Inductive Step: Assume the hypothesis for $k' < k$ and consider $k = H - height(I_i)$. Suppose that s_i is strictly dominated at I_i given $\Gamma_{-i}^{n-1}(I_i)$. Choose a strategy s'_i such that s'_i strictly dominates s_i given $\Gamma_{-i}^{n-1}(I_i)$ and s'_i is not strictly dominated at I_i given $\Gamma_{-i}^{n-1}(I_i)$. If s'_i is n -acceptable, the inductive claim is established. Otherwise let $I'_i > I_i$ be an information set such that (a) s'_i is strictly dominated at I'_i given $\Gamma_{-i}^{n-1}(I'_i)$ and (b) I'_i is of maximal height, i.e. if $I'_i > I_i^*$, then s'_i is not strictly dominated at I_i^* given $\Gamma_{-i}^{n-1}(I_i^*)$. By inductive hypothesis, we may choose a strategy s_i^* such that s_i^* strictly dominates s'_i given $\Gamma_{-i}^{n-1}(I'_i)$ and s_i^* is n -acceptable at I'_i . Applying the combination principle (Proposition 26) we may modify s'_i to obtain a strategy \widehat{s}_i such that \widehat{s}_i agrees with s_i^* on I'_i and agrees with s'_i outside of I'_i . From Lemma 17 it follows that \widehat{s}_i strictly dominates s_i at I_i given $\Gamma_{-i}^{n-1}(I_i)$. Finally, I note that if I_i^1, I_i^2, \dots are the maximal height information sets such that $I_i^j > I_i$ and s_i is strictly dominated at I_i^j , then we may successively construct \widehat{s}_i for I_i^1 , then \widehat{s}_i for I_i^2 , until we obtain a strategy \widehat{s}_i such that \widehat{s}_i strictly dominates s_i given $\Gamma_{-i}^{n-1}(I_i)$ and \widehat{s}_i is n -acceptable at I_i . This establishes the inductive claim.

To complete the proof of the proposition, let I_i^1, I_i^2, \dots be the information sets belonging to player i of maximal height, that is, for all I_i we have that $I_i^j \leq I_i$. Proceeding as in the inductive step, we may choose strategies s_i^j such that s_i^j is n -acceptable at I_i^j . Note that for each information set I_i , there is a unique predecessor I_i^j such that $I_i^j \leq I_i$. So we may define $s_i(I_i) = s_i^j(I_i)$ where I_i^j is the unique predecessor of I_i . Then s_i is n -acceptable at all information sets I_i in T , and hence $s_i \in \Gamma^n(T)$, as required. ■

Lemma 19 *Let T be a finite game tree with perfect recall. A strategy s_i is in $\Gamma_i^\infty(I) \iff$ for all information sets I_i such that s_i is consistent with I_i and I_i entails I for I :*

1. s_i is admissible at I_i given $S_{-i}(I_i)$, and
2. s_i is not strictly dominated at I_i given $\Gamma_{-i}^\infty(I_i)$.

Proof. (\Rightarrow) If s_i is not admissible at I_i given $S_{-i}(I_i)$, then s_i is not in $\Gamma_i^0(I)$ and hence not in $\Gamma_i^\infty(I)$. And if s_i is strictly dominated at I_i given $\Gamma_{-i}^\infty(I_i)$

then s_i must have been eliminated at a round before $\max(T)$, for otherwise the elimination procedure does not terminate at stage $\max(T)$.

(\Leftarrow) If s_i is admissible at all information sets I_i entailing I for i given $S_{-i}(T)$, and $s_i \in [I]_i$, then $s_i \in \Gamma_i^0(I)$. At later stages $n > 0$, suppose for reductio that $s_i \in [I]_i - \Gamma_i^n(I)$. Then s_i is strictly dominated at some information I_i given $\Gamma_i^{n-1}(I_i)$ where I_i entails I for i . Now by Proposition 18, we have $\Gamma_{-i}^\infty(I_i) \neq \emptyset$. Also, $\Gamma_{-i}^\infty(I_i) \subseteq \Gamma_i^{n-1}(I_i)$. So by Lemma 12 it follows that s_i is strictly dominated at I_i given $\Gamma_{-i}^\infty(I_i)$, which is a contradiction. Hence $s_i \in \Gamma_i^n(I_i)$ for all n , and thus $s_i \in \Gamma_i^\infty(I_i)$. ■

Lemma 20 *Let T be a game tree of perfect information with a unique subgame perfect equilibrium (s_i, s_{-i}) . Then a strategy profile (s'_i, s'_{-i}) is equal to $(s_i, s_{-i}) \iff$*

1. $s'_i(r)$ is a BI-maximizer, and
2. (s'_i, s'_{-i}) is an SPE in each subtree T_{r*a} .

Proof. (\Rightarrow) It is immediate that (s_i, s_{-i}) must be an SPE in each subtree. Now if $s_i(r)$ is not a BI-maximizer, consider the strategy s'_i that chooses a maximizer a such that $s'_i(r) = a$ and sets $s'_i(x) = s_i(x)$ for all $x \neq r$. Then by the second clause, (s'_i, s_{-i}) is an SPE in each subgame T_{r*a} . Since T has a unique subgame perfect equilibrium, this implies in particular that $u_i(s'_i, s_{-i}, r * a) = b_i(r * a)$. And clearly $u_i(s'_i, s_{-i}, r * a) = u_i(s'_i, s_{-i})$. By hypothesis, $b_i(r * a) > b_i(r * s_i(r)) = u_i(s_i, s_{-i})$. So s_i is not a best response to s_{-i} in T , contrary to the hypothesis that (s_i, s_{-i}) is in equilibrium.

(\Leftarrow) Since (s'_i, s'_{-i}) is an SPE in each subtree T_{r*a} , we have that $u_i(s'_i, s'_{-i}, r * a) = b_i(r * a)$ for all moves a at r . Since $s'_i(r)$ is a BI-maximizer, we have therefore that $u_i(s'_i, s'_{-i}, r * s'_i(r)) \geq u_i(s'_i, s'_{-i}, r * a)$ for all moves a at r . For every strategy s_i^* , we have that $u_i(s_i^*, s'_{-i}) = u_i(s_i^*, s'_{-i}, r * s_i^*(r)) \leq u_i(s'_i, s'_{-i}, r * s'_i(r))$, since (s'_i, s'_{-i}) is an equilibrium at $s'_i(r)$. All together, it follows that $u_i(s'_i, s'_{-i}) \geq u_i(s_i^*, s'_{-i})$. So s'_i is a best reply to s'_{-i} in T . Also s'_{-i} is a best reply to s'_i iff s'_{-i} is a best reply to s'_i in $T_{s'_i(r)}$, which is the case by Part 2. Hence (s'_i, s'_{-i}) form a Nash equilibrium in each subgame. Since by supposition (s_i, s_{-i}) is the only SPE in T , it follows that $(s_i, s_{-i}) = (s'_i, s'_{-i})$. ■

Proposition 21 *Let T be a finite game tree with perfect information and a unique subgame perfect equilibrium. Then for each node x , for each strategy profile $(s_i, s_{-i}) \in \Gamma^\infty(x)$, we have that*

1. (s_i, s_{-i}) is extendible to an SPE in T_x , and
2. $u_i(s_i, s_{-i}, x) = b_i(x)$.

Proof. The proof is by induction on the height h of node x . Part 1 immediately follows from Part 2.

Base Case, $h = 1$. Define $\max(i, x) = \max\{u_i(x * a) : x * a \text{ is in } T\}$. Suppose that *player*(x) = i . A move a for player i is a maximizer iff $x * a = \max(i, x)$.

Then s_i is strictly dominated at x given any subset of $S_{-i}(T)$ iff $s_i(x)$ is not a maximizer. Hence by Lemma 19, $(s_i, s_{-i}) \in \Gamma^\infty(x) \iff s_i(x)$ is a maximizer, and $s_i \in [x]_i$ and $s_{-i} \in [x]_{-i}$. This implies that each pair $(s_i, s_{-i}) \in \Gamma^\infty(x)$ is a subgame perfect equilibrium in T_x .

Inductive Step: Assume the hypothesis for h and consider $h + 1$. Let x be a node of height $h + 1$; by inductive hypothesis, we have for all successors $x * a$, for all profiles $(s_i, s_{-i}) \in \Gamma^\infty(x * a)$, that the payoff $u_i(s_i, s_{-i})$ equals $b_i(x * a)$. Note that all successors $x * a$ entail x both for i and $-i$. Hence it follows that (a) if $s_{-i} \in \Gamma^\infty(x)$, then $s_{-i} \in \Gamma^\infty(x * a)$ for all moves $x * a$ in T , and (b) if $(s_i, s_{-i}) \in \Gamma^\infty(x)$, then $(s_i, s_{-i}) \in \Gamma^\infty(x * s_i(x))$.

Let $s_i \in \Gamma_i^\infty(x)$. As before, define $\max(i, x) = \max\{b_i(x * a) : x * a \text{ is in } T\}$, and say that a move a for player i is a maximizer iff $x * a = \max(i, x)$. Suppose for reductio that $s_i(x)$ is not a maximizer, and let a be a maximizer. By Proposition 18, $\Gamma_i^\infty(x * a) \neq \emptyset$; let $s'_i \in \Gamma_i^\infty(x * a)$. I argue that s'_i strictly dominates s_i given $\Gamma_{-i}^\infty(x)$. Let $s_{-i} \in \Gamma_{-i}^\infty(x)$. Then by (a) $u_i(s'_i, s_{-i}, x) = u_i(s'_i, s_{-i}, x * s'_i(x)) = b_i(x * s'_i(x))$ by inductive hypothesis. And by (b), $(s_i, s_{-i}) \in \Gamma^\infty(x * s_i(x))$, and so $u_i(s_i, s_{-i}, x) = b_i(x * s_i(x))$. Since $s'_i(x)$ is a maximizer and $s_i(x)$ is not, we have that $u_i(s'_i, s_{-i}, x) > u_i(s_i, s_{-i}, x)$. As this holds for all $s_{-i} \in \Gamma_{-i}^\infty(x)$, it follows that s'_i strictly dominates s_i at x given $\Gamma_{-i}^\infty(x)$, and so by Lemma 19, $s_i \notin \Gamma_i^\infty(x)$, contrary to assumption. So if $s_i(x) \in \Gamma^\infty(x)$, then $s_i(x)$ is a maximizer.

Now let any profile $(s_i, s_{-i}) \in \Gamma^\infty(x)$ be given. Choose a strategy s'_i such that for each $a \neq s_i(x)$ we have that $(s'_i, s_{-i}) \in \Gamma^\infty(x * a)$. Define $s_i^*(x') = s_i(x)$ if $x' > x * s_i(x)$, and $s_i^*(x') = s'_i(x')$ otherwise. Then $(s_i^*, s_{-i}) \in \Gamma^\infty(x * a)$ for all moves a at x ; hence by inductive hypothesis we may choose strategies $(\widehat{s}_i, \widehat{s}_{-i})$ such that in each subtree T_{x*a} : \widehat{s}_i^* extends s_i^* and \widehat{s}_{-i} extends s_{-i} and $(\widehat{s}_i, \widehat{s}_{-i})$ is an SPE in T_{x*a} . So provided that $\widehat{s}_i(x) = s_i(x)$, it follows from Lemma 2 that $(\widehat{s}_i, \widehat{s}_{-i})$ is an SPE in T_x . So $(\widehat{s}_i, \widehat{s}_{-i})$ satisfies the requirements of the Proposition. This concludes the inductive step and establishes the claim. \blacksquare

Lemma 24 *Let T be a game tree with perfect recall. Then*

1. *for all nodes x, x' if $x < x'$ then $I(x) \neq I(x')$.*
2. *if $s_i \in [I_i]_i$ and $(s'_i, s_{-i}) \in I_i$, then $I_i \cap \text{Play}(s_i, s_{-i}) = I_i \cap \text{Play}(s'_i, s_{-i})$.*
3. *if $(s_i, s_{-i}) \in [I_i]$ and $(s'_i, s'_{-i}) \in [I_i]$, then $(s_i, s'_{-i}) \in [I_i]$.*

Proof. Part 1. Suppose otherwise. Then let x, x' be such that $x < x'$ and $I(x) = I(x')$. Let $a = \text{move}(x, x')$. Then $E_i(x')$ extends $E_i(x) * a$. Hence $E_i(x') \neq E_i(x)$, contrary to the hypothesis that T is a game tree with perfect recall.

Part 2. From part 1 it follows that $\text{Play}(s'_i, s_{-i})$ intersects I_i in exactly one node; let $I_i \cap \text{Play}(s'_i, s_{-i}) = \{x\}$. Let $s_i \in [I_i]_i$; then there is a strategy s'_{-i} such that $\text{play}(s_i, s'_{-i})$ reaches I_i . Let $y \in \text{Play}(s_i, s'_{-i}) \cap I_i$. I now argue by induction that $\text{Play}(s_i, s_{-i})$ contains all prefixes of x . Base Case: Clearly

$\emptyset \in \text{Play}(s_i, s_{-i})$. Inductive Step: Suppose that $x' < x \in \text{Play}(s_i, s_{-i})$; let $a = \text{move}(x, x')$. Case 1: $\text{player}(x') = -i$. Then $a = s_{-i}(x')$, and so $x' * a \in \text{Play}(s_i, s_{-i})$. Case 2: $\text{player}(x') = i$. Then $E_i(x)$ is of the form $E_i(x') * I(x') * s'_i(x') * \dots$. By perfect recall, $E_i(x) = E_i(y)$ and so $E_i(y)$ begins with $E_i(x') * I(x') * s'_i(x') * \dots$. Since $y \in \text{Play}(s_i, s'_{-i})$, it follows that $s_i(x') = s'_i(x') = a$. Hence $x' * a \in \text{Play}(s_i, s'_{-i})$. This shows that all prefixes of x , including x itself, are in $\text{Play}(s_i, s_{-i})$.

Part 3. Suppose that $(s_i, s_{-i}) \in [I_i]$ and $(s'_i, s'_{-i}) \in [I'_i]$. Then by Part 2, $I_i \cap \text{Play}(s_i, s'_{-i}) = I_i \cap \text{Play}(s'_i, s'_{-i})$. Since $\text{Play}(s'_i, s'_{-i}) \cap I_i \neq \emptyset$, we have that $I_i \cap \text{Play}(s_i, s'_{-i}) \neq \emptyset$; hence $(s_i, s_{-i}) \in [I'_i]$, as required. ■

Lemma 25 *Let T be a game tree with perfect recall. Let I_i be an information set belonging to player i . Then*

1. for all $s_i \in [I_i]_i$ we have that $\text{cons}(s_i, I_i) = [I_i]_{-i}$
2. for all information sets I'_i , if $I_i \geq I'_i$, then $[I_i] \subseteq [I'_i]$
3. for all information sets I , if $I_i \geq I$ then $[I_i]_i \subseteq [I]_i$
4. for all information sets I_i, I'_i of player i it is the case that: $I_i = I'_i \iff (I_i \geq I'_i \text{ and } I'_i \geq I_i)$

Proof. Part 1. Let $s_i \in [I_i]_i$ be a strategy consistent with I_i ; then there is a strategy s_{-i} such that $(s_i, s_{-i}) \in [I_i]$. Let s'_{-i} be any strategy for player $-i$ consistent with I_i ; then there is a strategy s'_i such that $(s'_i, s'_{-i}) \in [I_i]$. Hence by Part 3 of Lemma 24, $(s_i, s_{-i}) \in [I_i]$. Since s_{-i} was chosen arbitrarily, it follows that $\text{cons}(s_i, I_i) = [I_i]_{-i}$.

Part 2. If $I_i = I'_i$, the claim follows immediately. Otherwise choose nodes $x \in I_i, x' \in I'_i$ such that $x > x'$. Let $(s_i, s_{-i}) \in [I_i]$ and let y be some node in $\text{Play}(s_i, s_{-i}) \cap I_i$. Then by perfect recall $E_i(y) = E_i(x)$, so I'_i appears in $E_i(y)$. Hence there is a node $y' \in I'_i$ such that $y' < y$, which implies that $y' \in \text{play}(s_i, s_{-i})$. So $(s_i, s_{-i}) \in [I'_i]$, which establishes that $[I_i] \subseteq [I'_i]$.

Part 3. If $I_i = I$ it follows immediately that $[I_i]_i = [I]_i$. So suppose that $I_i > I$. Let $x > x'$ where $x \in I_i, x' \in I$. Let $s_i \in [I_i]_i$. By the basic reachability assumption (cf. Section 2), there is a strategy pair (s'_i, s_{-i}) such that $x \in \text{Play}(s'_i, s_{-i})$. By Part 1 of the current lemma, it follows that $s'_{-i} \in \text{cons}(s_i, I_i)$ and hence $(s_i, s_{-i}) \in [I_i]$. By Part 2 of Lemma 24, it follows that $I_i \cap \text{Play}(s_i, s_{-i}) = I_i \cap \text{Play}(s'_i, s_{-i})$. Hence $x \in \text{Play}(s_i, s_{-i})$, and therefore $x' \in \text{Play}(s_i, s_{-i})$. So $s_i \in [I(x')]_i = [I]_i$.

Part 4. (\implies) Immediate. (\impliedby) Suppose that $I_i > I'_i$ and $I'_i > I_i$ but $I_i \neq I'_i$. So there is $x, y \in I_i$ and $x', y' \in I'_i$ such that $x > x'$ and $y' > y$. Then $I(x) = I_i$ appears in the epistemic history $E_i(x')$. Since $I(x') = I(y') = I'_i$, by perfect recall it follows that $E_i(x') = E_i(y')$. Hence I_i appears in $E_i(y')$ and thus in $E_i(y)$ since $y' > y$ and hence $E_i(y)$ is of the form $E_i(y') * \dots$. Thus there is a node $y_0 > y' > y$ such that $I(y_0) = I(y) = I_i$. By Part 1 of Lemma 24, this contradicts perfect recall. ■

Proposition 26 (Combination Principle) *Let T be a game tree with perfect recall. Let s_i, s'_i be two strategies consistent with an information set I_i (i.e., $s_i, s'_i \in [I_i]_i$). Then there is a strategy s_i^* such that*

1. s_i^* agrees with s'_i on I_i (i.e., for all strategies s_{-i} if $(s_i, s_{-i}) \in [I_i]$, then $\text{play}(s_i^*, s_{-i}) = \text{play}(s'_i, s_{-i})$), and
2. s_i^* agrees with s_i outside of I_i (i.e., for all strategies s_{-i} , if $(s_i, s_{-i}) \notin [I_i]$, then $\text{play}(s_i^*, s_{-i}) = \text{play}(s_i, s_{-i})$).

Proof. Define s_i^* :

$$s_i^*(I_i) = \begin{cases} s'_i(I_i) & \text{if } I_i \geq I_i \\ s_i(I_i) & \text{o.w.} \end{cases} .$$

Part 1. Suppose that $(s_i, s_{-i}) \in [I_i]$. I show that $\text{play}(s_i^*, s_{-i}) = \text{play}(s'_i, s_{-i})$. By Part 1 of Lemma 24, $\text{play}(s_i, s_{-i}) \cap I = \{x\}$ for some node $x \in I$. So from Part 2 of Lemma 24, we have that $(*) \text{Play}(s'_i, s_{-i}) \cap I = \text{Play}(s_i, s_{-i}) \cap I = \{x\}$. I show that $x \in \text{Play}(s_i^*, s_{-i})$. Let $x' < x$ such that $\text{player}(x') = i$. Then $I(x') \neq I(x)$ by Part 1 of Lemma 24. Also, $I(x) \geq I(x')$. So by Part 4 of Lemma 25, $I(x') \not\geq I(x) = I_i$. Hence from the definition of s_i^* , it follows that $s_i^*(I(x')) = s_i(I(x'))$. Since this holds for all $x' < x$, it follows that $x \in \text{Play}(s_i^*, s_{-i})$. Given $(*)$, this shows that $\text{play}(s_i^*, s_{-i})$ agrees with $\text{play}(s'_i, s_{-i})$ up to x .

I now show that $x' \in \text{Play}(s_i^*, s_{-i})$ for all $x' > x$ such that $x' \in \text{Play}(s'_i, s_{-i})$.

Base Case, $x' = x * a$. Since $I(x) = I_i$, it follows from the definition of s_i^* that $s_i^*(I(x)) = s'_i(I(x)) = a$ and hence $x' \in \text{Play}(s_i^*, s_{-i})$.

Inductive Step: Assume the hypothesis for x^* and consider $x' = x^* * a$.

Case 1: $\text{player}(x^*) \neq i$. Then $s_{-i}(x^*) = a$, and so $x^* * a \in \text{Play}(s'_i, s_{-i})$.

Case 2: $\text{player}(x^*) = i$. Since $x^* \geq x$, it follows that $I(x^*) \geq I(x) = I_i$. So by definition of s_i^* , we have that $s_i^*(I(x^*)) = s'_i(I(x^*)) = a$. Hence $x^* * a \in \text{Play}(s_i^*, s_{-i})$. This completes the inductive step and establishes Part 1.

For Part 2, I show the contrapositive. Suppose that $\text{play}(s_i, s_{-i}) \neq \text{play}(s_i^*, s_{-i})$ for some strategy s_{-i} . Since $\emptyset \in \text{Play}(s_i, s_{-i}) \cap \text{Play}(s_i^*, s_{-i})$, there is a greatest x such that $x \in \text{Play}(s_i, s_{-i}) \cap \text{Play}(s_i^*, s_{-i})$. (The sequences $\text{play}(s_i, s_{-i})$ and $\text{play}(s_i^*, s_{-i})$ agree up to x and then diverge.) Then $\text{player}(x) = i$ for otherwise s_{-i} would choose the same action at x against both s_i and s_i^* , that is, $x * s_{-i}(x) \in \text{Play}(s_i, s_{-i}) \cap \text{Play}(s_i^*, s_{-i})$. Similarly, $s_i(I(x)) \neq s_i^*(I(x))$. So by definition of s_i^* , it follows that $I(x) \geq I_i$. So by Part 2 of Lemma 25, it follows that $[I(x)] \subseteq [I_i]$. Since $x \in \text{play}(s_i, s_{-i})$, it follows that $(s_i, s_{-i}) \in [I(x)]$ and thus $(s_i, s_{-i}) \in [I_i]$. So if $\text{play}(s_i, s_{-i}) \neq \text{play}(s_i^*, s_{-i})$, then $(s_i, s_{-i}) \in [I_i]$. Hence s_i^* agrees with s_i outside of I_i , which was to be shown. ■

Proposition 27 *Let T be a game tree with perfect recall. Let ρ_i be an LPS for player i (i.e., ρ_i defines the preferences of player i over strategies) with full support. Consider any information set I_i . Suppose that (1) s_i is weakly dominated at I_i , or (2) s_i is strictly dominated at I_i given Σ_{-i} and $\rho_i * [I_i]_{-i} \subseteq \Sigma_{-i}$. Then there is a strategy s'_i such that*

1. $s'_i \in [I_i]_i$ and
2. $s'_i \sim^{[I_i]_{-i}} s_i$, and
3. $s'_i \succ^{[I_i]_{-i}} s_i$, and
4. $s'_i \succ_i s_i$.

Proof. Part 1. Suppose that s_i^* weakly dominates s_i at I_i or strictly dominated s_i at I_i given Σ_{-i} . Then Definition 11 implies that $s_i, s_i^* \in [I_i]_i$. So by Proposition 26, there is a strategy s'_i that agrees with s_i^* on I_i and agrees with s_i outside of I_i . Since s'_i agrees with s_i^* on I_i , it follows that $s'_i \in [I_i]_i$, which establishes Part 1.

Part 4 follows from Parts 2, 3 and Proposition 4.

To establish Part 2, let $s_{-i} \in \overline{[I_i]_{-i}}$ be given. By Clause 1 of Lemma 24, $[I_i]_{-i} = \text{cons}(s_i, s_{-i})$ and so $(s_i, s_{-i}) \notin [I_i]$. Since s'_i agrees with s_i outside of I_i , it follows that $\text{play}(s'_i, s_{-i}) = \text{play}(s_i, s_{-i})$. Hence for all $s_{-i} \in \overline{[I_i]_{-i}}$ we have $u_i(s_i, s_{-i}) = u_i(s'_i, s_{-i})$, so $s'_i \sim^{[I_i]_{-i}} s_i$ as required.

For Part 3, I consider two cases.

Case 1: s_i^* weakly dominates s_i at I_i . Let $s_{-i} \in [I_i]_{-i}$. Then by Clause 1 of Lemma 24, $(s_i^*, s_{-i}) \in [I_i]$. Since s'_i agrees with s_i^* on I_i , it follows that $u_i(s'_i, s_{-i}) = u_i(s_i^*, s_{-i})$. Moreover, $u_i(s_i^*, s_{-i}) \geq u_i(s_i, s_{-i})$ since s_i^* weakly dominates s_i at I_i . For the same reason, there is some strategy $s'_{-i} \in [I_i]_{-i}$ such that $u_i(s_i^*, s'_{-i}) > u_i(s_i, s'_{-i})$. As before we have that $u_i(s_i^*, s'_{-i}) = u_i(s'_i, s'_{-i})$, so s'_i weakly dominates s_i in the space reduced to $[I_i]_{-i}$. It therefore follows from Lemma 2 that $s'_i \succ_i^{[I_i]} s_i$.

Case 2: s_i^* strictly dominates s_i at I_i given Σ_{-i} . Then for any $s_{-i} \in \Sigma_{-i}$ such that $(s_i, s_{-i}) \in [I_i]$, we have that $u_i(s'_i, s_{-i}) = u_i(s_i^*, s_{-i}) = u_i(s_i, s_{-i})$. Hence s'_i strictly dominates s_i at I_i given Σ_{-i} . Now suppose that $\rho_i * [I_i]_{-i} \subseteq \Sigma_{-i}$. Since ρ_i has full support, it follows that $\rho_i * [I_i]_{-i} \neq \emptyset$, and so by Lemma 12, we have that (*) s_i^* strictly dominates s_i given $\rho_i * [I_i]_{-i}$, which is the support of $(\rho_i | [I_i]_{-i})^1$. So (*) implies that $EU((\rho | [I_i]_{-i})^1, s'_i, u_i) > EU((\rho | [I_i]_{-i})^1, s_i, u_i)$, which implies that $s'_i \succ_i s_i$.

Hence in either case $s'_i \succ_i s_i$, which establishes the Proposition. ■

References

- [1] Geir Asheim. Proper rationalizability in lexicographic beliefs. *International Journal of Game Theory*, 30:453:478, 2001.
- [2] B.D. Bernheim. Rationalizable strategic behavior. *Econometrica*, 52:1007–1028, 1984.
- [3] Lawrence Blume, Adam Brandenburger, and Eddie Dekel. Lexicographic probabilities and choice under uncertainty. *Econometrica*, 59(1):61–79, 1991.

- [4] Lawrence Blume, Adam Brandenburger, and Eddie Dekel. Lexicographic probabilities and equilibrium refinements. *Econometrica*, 59(1):81–98, 1991.
- [5] Adam Brandenburger and Eddie Dekel. Hierarchies of beliefs and common knowledge. *Journal of Economic Theory*, 59:189–198, 1993.
- [6] Eddie Dekel and Drew Fudenberg. Rational behavior with payoff uncertainty. *Journal of Economic Theory*, 52:243–267, 1990.
- [7] Peter Gärdenfors. *Knowledge In Flux: Modeling the Dynamics of Epistemic States*. MIT Press, Cambridge, Mass., 1988.
- [8] Mamoru Kaneko. Epistemic logics and their game theoretic applications: Introduction. *Economic Theory*, 19:1:7–62, 2002.
- [9] Mamoru Kaneko and J. Jude Kline. Behavior strategies, mixed strategies and perfect recall. *International Journal of Game Theory*, 24:127–145, 1995.
- [10] Elon Kohlberg. Refinement of nash equilibrium: The main ideas. In *Game Theory and Applications*. Academic Press, San Diego, 1990.
- [11] J-M Mertens and S. Zamir. 1985. *Formulation of Bayesian analysis for games of incomplete information*, 14:1–29, 1985.
- [12] R. Myerson. Refinements of the nash equilibrium concept. *International Journal of Game Theory*, 7:73–80, 1978.
- [13] Martin J. Osborne and Ariel Rubinstein. *A Course in Game Theory*. MIT Press, Cambridge, Mass., 1994.
- [14] D.G. Pearce. Rationalizable strategic behavior and the problem of perfection. *Econometrica*, 52(1029–1050), 1984.
- [15] Hans Rott. Coherence and conservatism in the dynamics of belief. part i: Finding the right framework. *Erkenntnis*, 50:387–412, 1999.
- [16] F. Schuhmacher. Proper rationalizability and backward induction. *International Journal of Game Theory*, 28:599–615, 1999.
- [17] Oliver Schulte. Minimal belief change, pareto-optimality and logical consequence. *Economic Theory*, 19(1):105–144, 2002.
- [18] Robert Stalnaker. Knowledge, belief and counterfactual reasoning in games. *Economics and Philosophy*, 12:133–163, 1996.
- [19] Robert Stalnaker. On the evaluation of solution concepts. *Theory and Decision*, 37:49–73, 1996.
- [20] Robert Stalnaker. Belief revision in games: forward and backward induction. *Mathematical Social Sciences*, 36:31–56, 1998.