

Lexicographic agreeing to disagree and perfect equilibrium<sup>☆</sup>Christian W. Bach<sup>a,b,\*</sup>, Jérémie Cabessa<sup>c</sup><sup>a</sup> Department of Economics, University of Liverpool Management School, University of Liverpool, Chatham Street, Liverpool, L69 7ZH, United Kingdom<sup>b</sup> EPICENTER, School of Business and Economics, Maastricht University, 6200 MD Maastricht, The Netherlands<sup>c</sup> DAVID Laboratory, University of Versailles Saint-Quentin-en-Yvelines, University of Paris-Saclay, 45 Avenue des États-Unis, 78035 Versailles, France

## ARTICLE INFO

Manuscript handled by Editor Andrés Carvajal

## Keywords:

Agreeing to disagree  
Lexicographic probability systems  
Epistemic game theory  
Mutual absolute continuity  
Perfect equilibrium  
Static games

## ABSTRACT

Aumann's seminal agreement theorem deals with the impossibility for agents to acknowledge their distinct posterior beliefs. We consider agreeing to disagree in an extended framework with lexicographic probability systems. A weak agreement theorem in the sense of identical posteriors only at the first lexicographic level obtains. Somewhat surprisingly, a possibility result does emerge for the deeper levels. Agents can agree to disagree on their posteriors beyond the first lexicographic level. By means of mutual absolute continuity as an additional assumption, a strong agreement theorem with equal posteriors at every lexicographic level ensues. Subsequently, we turn to games and provide epistemic conditions for the classical solution concept of perfect equilibrium. Our lexicographic agreement theorems turn out to be pivotal in this endeavour. The hypotheses of mutual primary belief in caution, mutual primary belief in rationality, and common knowledge of conjectures characterize perfect equilibrium epistemically in our lexicographic framework.

## 1. Introduction

The impossibility for two agents to agree to disagree is established by Aumann (1976)'s seminal agreement theorem. More precisely, if two Bayesian agents with a common prior receive private information and have common knowledge of their posterior beliefs, then these posteriors must be equal. In other words, distinct posterior beliefs cannot be common knowledge among Bayesian agents with the same prior beliefs. In this sense, agents cannot agree to disagree.<sup>1</sup> The impossibility of agreeing to disagree has important implications for any interactive situation where, loosely speaking, the mutual acknowledgement of distinct views or assessments is relevant, e.g. trade, speculation, political

positions, or legal judgements.<sup>2</sup> The array of potential applications for the agreement theorem is vast.

Here, we explore agreeing to disagree in an extended framework with lexicographic beliefs. A lexicographic belief is a sequence of beliefs, where the different beliefs are given in descending order of importance.<sup>3</sup> The sequence's first component can be viewed as the agent's primary doxastic attitude, its second component as his secondary doxastic attitude, etc. Intuitively, a lexicographically-minded agent deems his first belief fundamentally more likely than his secondary belief, which in turn is fundamentally more likely than his tertiary belief, etc. Lexicographic beliefs resolve the problem of conditioning on events with probability zero. Revising beliefs based on

<sup>☆</sup> Preliminary versions of this work were presented at the 13th Conference on Logic and the Foundations of Game and Decision Theory (LOFT13), Milan, July 2018, as well as at the 6th World Congress of the Game Theory Society (GAMES2020), Budapest, July 2021. We are grateful to Adam Brandenburger, Andrés Carvajal, Robert Edwards, Amanda Friedenberg, Stephan Jagau, Andrés Perea, Burkhard Schipper, Elias Tsakas, and two anonymous referees for useful as well as constructive comments.

\* Corresponding author at: Department of Economics, University of Liverpool Management School, University of Liverpool, Chatham Street, Liverpool, L69 7ZH, United Kingdom.

E-mail address: [cwbach@liverpool.ac.uk](mailto:cwbach@liverpool.ac.uk) (C.W. Bach).

<sup>1</sup> An extensive literature on agreeing to disagree has emerged. Most contributions reconsider Aumann's impossibility theorem in more general frameworks. Notably, Bonanno and Nehring (1997) as well as Ménager (2012) provide comprehensive surveys on this literature. Some more recent contributions to the agreeing to disagree literature include Dégremon and Roy (2012), Hellman and Samet (2012), Bach and Perea (2013), Heifetz et al. (2013), Hellman (2013), Demey (2014), Lehrer and Samet (2014), Chen et al. (2015), Dominiak and Lefort (2015), Tarbush (2016), Bach and Cabessa (2017), Gizatulina and Hellman (2019), Pacuit (2018), Tsakas (2018), Liu (2019), as well as Contreras-Tejada et al. (2021).

<sup>2</sup> A prominent analysis of economic consequences of agreeing to disagree is Milgrom and Stokey's (1982) so-called no-trade theorem. Accordingly, if two traders agree on a prior efficient allocation of goods, then upon receiving private information it cannot be common knowledge that they both have an incentive to trade.

<sup>3</sup> Formally, lexicographic beliefs are modelled in their most general form by lexicographic probability systems due to Blume et al. (1991a).

	<i>y</i>	<i>z</i>
<i>a</i>	1, 0	0, 1
<i>b</i>	0, 0	0, 1

Fig. 1. A two player game.

hypotheses that are initially deemed impossible is relevant to hypothetical reasoning. An apt example are games. It can be important for a player to consider what would happen, if an opponent were to pick an unexpected choice, in order to act rationally himself.

In game theory, lexicographic beliefs do play a prominent role and have effectively been put into action to model caution and trembles.<sup>4</sup> In particular, they shed essential light on the foundations of weak dominance arguments and have served to unravel a fundamental game-theoretic paradox: the so-called inclusion-exclusion problem.<sup>5</sup> The paradox arises whenever a player is required to *include* all, yet to *exclude* some, choices for an opponent. This startling tension is inherent in (iterated) weak dominance, also called (iterated) admissibility, which constitutes one of the most long-standing ideas in game theory going back at least to Gale (1953).

For an illustration of the inclusion-exclusion problem, consider the two player game depicted in Fig. 1 with players *Alice* and *Bob*, where *Alice* chooses a “row” (*a* or *b*) and *Bob* picks a “column” (*y* or *z*). The unique strategy for *Alice* in line with weak dominance is *a*. Intuitively, against all choices of *Bob*, *a* never yields less than *b*, and against the particular strategy *y* of *Bob*, *a* induces a strictly higher payoff than *b*. For *Bob*, *y* is strictly worse than *z* against all of *Alice*’s choices. However, it seems impossible to support *a* with consistent beliefs, since on the one hand, *Alice* needs to assign positive probability to both *y* and *z* to render *a* uniquely optimal for her, while on the other hand, she should assign probability zero to the never optimal choice *y* for *Bob*. The remedy to the paradox lies in lexicographic beliefs. They are capable of not excluding any choice from consideration yet at the same time deeming some choices much more – indeed infinitely more – likely than others. With lexicographic beliefs, the inclusion-exclusion riddle evaporates. In the preceding example, a lexicographic belief for *Alice* that assigns probability one to *z* in its first level and probability one to *y* in its second level would already form a consistent doxastic attitude filtering out *a* as her unique optimal strategy.

In terms of Aumann’s impossibility theorem the question of whether agreeing to disagree is possible or not gains in depth if lexicographic beliefs are admitted and hypothetical reasoning can thereby be captured. For example, consider merchants forming beliefs about the arrival of a sea shipment. A primary contingency could revolve around the usual meteorological conditions that can affect the length of sea travel. Suppose that a secondary contingency would include fundamentally less likely factors affecting arrival like a pirate attack. If common knowledge of their posterior beliefs implies agents to agree on their beliefs given the primary contingency, then they could possibly still disagree with regards to the secondary contingency. Whether or not the agents do, could have different implications for the actions they take based on their (lexicographic) beliefs.

In general, given the importance of lexicographic beliefs in game theory on the one hand, and given Aumann’s seminal impossibility

result on agreeing to disagree on the other hand, it seems intriguing to ask how the agreement theorem is affected if standard probabilities are replaced by lexicographic probability systems. To address this question we define the notion of lexicographic Aumann structure, where the agents hold a sequence of priors on the basis of which they compute a sequence of posteriors in the style of Blume et al. (1991a). In our framework, a weak agreement theorem in the sense of merely identical first level posteriors obtains. However, we provide a disagreement result establishing that agents can actually agree to disagree on their posteriors beyond the first lexicographic level. Aumann’s impossibility theorem does therefore not directly generalize to full-fledged lexicographic reasoning. Based on this observation, we introduce a condition which essentially states that every lexicographic level prior either neglects or considers the agents’ private information synchronically. This condition can be viewed as a variant of standard mutual absolute continuity from probability theory. With the assistance of mutual absolute continuity, we provide a strong agreement theorem which establishes the impossibility of agreeing to lexicographically disagree.

Naturally, the question arises whether our lexicographic agreement theorems can be applied to game theory. It would be particularly illuminating to gain novel insights about classical solution concepts based on lexicographic agreeing to disagree. A prominent class of solution concepts in game theory is based on the idea of trembles. Intuitively, with a very small probability a player may make a mistake – “his hand might tremble” – in implementing his optimal strategy. So-called tremble equilibria formalize this intuition by postulating equilibrium behaviour as the limiting case when the trembles vanish. The most fundamental solution concept of this kind is Selten’s (1975) perfect equilibrium.<sup>6</sup> A typical feature of tremble equilibria requires all trembles to satisfy some full support condition. In this sense, tremble equilibria also formalize cautious players, which suggests a link to lexicographic beliefs. Indeed, Blume et al. (1991b) investigate this link and provide a reformulation of perfect equilibrium as well as of proper equilibrium in terms of lexicographic conjectures, which are lexicographic beliefs about choices.

However, a characterization of tremble equilibria in terms of interactive thinking is still missing. Such an endeavour would imperatively involve higher-order beliefs, thereby moving beyond the basic doxastic layer of conjectures. Full interactive reasoning is modelled by imposing conditions on belief hierarchies which in turn assemble different layers of iterated beliefs. Conjectures, as beliefs about (opponents’) choices, only constitute the first such layer. In order to fully describe the interactive thinking of players, it is crucial to also model their beliefs about their opponents’ conjectures, their beliefs about their opponents’ beliefs about their opponents’ conjectures, etc. Due to their infinite nature belief hierarchies are cumbersome objects, but fortunately they can be represented in a compact way by means of epistemic models due to Harsanyi (1967–68). The epistemic programme in game theory has employed such models to unveil the interactive reasoning assumptions implicitly endorsed by solution concepts in games.

Our lexicographic agreement theorems are capable of shedding some light on the interactive reasoning underlying perfect equilibrium in games. Indeed, we provide epistemic conditions for perfect equilibrium. The epistemic hypotheses of mutual primary belief in caution, mutual primary belief in rationality, and common knowledge of conjectures characterize perfect equilibrium in terms of interactive reasoning. Our lexicographic agreement theorems play a prominent role in attaining our epistemic foundation. By means of the weak agreement theorem, all opponents of any given player can be ensured to

<sup>4</sup> By now lexicographic beliefs have become a widespread tool in game theory and have been used, for instance, by Kreps and Wilson (1982), Kreps and Ramey (1987), Blume et al. (1991b), Brandenburger (1992b), Börgers (1994), Stahl (1995), Mailath et al. (1997), Asheim (2001, 2002), Govindan and Klumpp (2003), Asheim and Perea (2005), Brandenburger et al. (2008), Yang (2015), Dekel et al. (2016), Lee (2016), as well as Catonini and De Vito (2018, 2020).

<sup>5</sup> The inclusion-exclusion problem has first been identified by Samuelson (1992), when showing that the solution concept of iterated weak dominance can be inconsistent with common knowledge assumptions.

<sup>6</sup> Other tremble equilibria have been proposed in the literature, for instance, Myerson’s (1978) proper equilibrium, van Damme’s (1984) quasi-perfect equilibrium, as well as Harsanyi and Selten’s (1988) uniformly perfect equilibrium.

hold the same marginal lexicographic conjecture about him. The strong agreement theorem is used to derive an independence property of the players' lexicographic conjectures.

We proceed as follows. The remainder of this section demarcates our model and results from the related literature. In Section 2, Blume et al.'s (1991a) lexicographic probability systems are incorporated into state-based interactive epistemology. Core notation is fixed and key concepts are defined. Section 3 contains a weak agreement theorem (WAT) with lexicographic probability systems, while Section 4 brings the deeper lexicographic levels into focus. Incongruity can obtain beyond the first level as our disagreement result (DIS) shows. In Section 5, under mutual absolute continuity, a lexicographically strong agreement theorem (SAT) is developed. We subsequently turn to games. In Section 6, Selten's (1975) seminal solution concept of perfect equilibrium is presented. A reformulation of this tremble equilibrium by means of lexicographic conjectures is furnished along the lines of Blume et al. (1991b) in Section 7. Epistemic conditions that characterize perfect equilibrium are put forth in Section 8. Finally, Section 9 offers some concluding remarks.

### 1.1. Related literature

By establishing agreement theorems with lexicographic beliefs and providing epistemic conditions for perfect equilibrium, our contribution is twofold. On the one hand, we are connected to the literature on agreeing to disagree that has emerged since Aumann's seminal (1976) impossibility result. On the other hand, the application of our lexicographic agreement theorems to epistemically characterize perfect equilibrium adds to the foundations of game theory.

Our framework extends standard Aumann structures (Aumann, 1974, 1976) by modelling the agents' beliefs with Blume et al.'s (1991a) lexicographic probability systems instead of mere probability distributions. Within this enriched set-up, we explore agreeing to disagree. Aumann's (1976) agreement theorem obtains as a special case of WAT, if the lexicographic common prior is truncated at the first level.

A lexicographic approach to agreeing to disagree is also taken by Bach and Perea (2013). Notably, their framework admits lexicographic beliefs as priors yet delivers a standard posterior for every agent. In contrast, by using lexicographic probability systems, we also model the posteriors as lexicographic beliefs. This does not only formally but also conceptually make an essential difference, as the agents' decision-relevant beliefs are the posteriors which are extended in our framework. A further restriction of Bach and Perea (2013) is a non-overlapping support requirement on lexicographic priors, which we do not impose. The agreement theorem of Bach and Perea (2013) is implied as another special case of WAT, if the lexicographic posteriors are truncated at the first level.

Once lexicographic posteriors enter the picture novel insights emerge. Somewhat surprisingly, our possibility result DIS establishes that agents can actually agree to disagree with a lexicographic mindset. In fact, if a non-overlapping support requirement on lexicographic priors were to be desired, DIS would still remain valid. The additional assumption of mutual absolute continuity brings about our impossibility result SAT, which can be viewed as a *lexicographic* agreement theorem in *sensu stricto*.

In general, lexicographic probability systems deal with the problem of how to proceed if something is learned to which initially probability zero was assigned. An alternative tool for extending probabilities to handle conditioning on measure zero events are conditional probability systems due to Rényi (1995). They have prominently been used in game theory to define the reasoning concept of common strong belief in rationality for extensive forms by Battigalli and Siniscalchi (2002). Lexicographic probability systems can be related to conditional probability systems and equivalences have been established under certain conditions (e.g. Hammond, 1994; Halpern, 2010; Tsakas, 2014). Lexicographic agreeing to disagree is thus indirectly also related to Tsakas

(2018), who establishes two agreement theorems with conditional probability systems. However, his results cannot be directly compared to ours, since the models are too different. While we extend Aumann's partitional model by lexicographic probability systems, Tsakas (2018) uses type structures in the style of Battigalli and Siniscalchi (1999). In particular, the way in which the agents' posteriors enter the picture is inherently distinct. In Tsakas (2018) framework, the agreement concerns a single posterior per agent, while our agreement theorems deal with lexicographic posteriors. Besides, already the computation of the first level posterior in our framework depends on which prior assigns positive probability to the conditioning event (i.e. the respective agent's information cell in lexicographic Aumann structures). In contrast, the determination of the conditioning event to derive the posterior in Tsakas (2018) model is independent from the prior.

In the game-theoretic part of our paper, we explore the epistemic foundation of Selten's (1975) solution concept of perfect equilibrium. A reformulation of perfect equilibrium by means of lexicographic conjectures constitutes the first step. Although such a reformulation has already been established by Blume et al. (1991b), our Lemma 1 provides a similar construction for the sake of completeness and self-containedness. Being concerned with the players' interactive reasoning, epistemic foundations go beyond conjectures into the players' belief hierarchies. Our Theorems 3 and 4 provide an epistemic characterization of perfect equilibrium. They can be viewed as developing Blume et al.'s (1991a) analysis of perfect equilibrium in terms of lexicographic conjectures further into the full game-theoretic reasoning realm. In some sense, our relation to Blume et al. (1991b) with regard to perfect equilibrium is analogous to the relation of Aumann and Brandenburger (1995) to Harsanyi (1973) with regard to Nash equilibrium: while Harsanyi (1973) has proposed the interpretation of Nash equilibrium in terms of conjectures, Aumann and Brandenburger (1995) have taken this crucial insight into an epistemic framework, unveiling the underlying interactive reasoning assumptions of Nash equilibrium. Our game-theoretic results could be perceived of as generalizing Aumann and Brandenburger (1995) from Nash equilibrium to perfect equilibrium.<sup>7</sup>

For the special case of two players, perfect equilibrium has been characterized epistemically by Perea (2012). The supply of epistemic conditions for perfect equilibrium involving any finite number of players has still been an open question though, which our Theorems 3 and 4 address. An epistemic analysis of equilibrium notions faces two considerable challenges once more than two players are considered. Firstly, for a given player, all opponents have to share the same belief about the player's choice ("problem of projection"). Secondly, any player's belief about his opponents' choices needs to be independent ("problem of independence"). Our lexicographic agreement theorems turn out to be pivotal in resolving these intricacies. Besides his restriction to the two player case, Perea's (2012) type-based framework is distinct from our state-based lexicographic Aumann structures with a common prior. Epistemic conditions for the special setting of two players are provided by our Proposition 2, which can thus be juxtaposed with Perea (2012). Our hypotheses of mutual primary belief in caution and mutual primary belief in rationality are weaker variants of his common full belief in caution and common full belief in primary belief in rationality, respectively. Furthermore, mutual knowledge of lexicographic conjectures embodies a correct beliefs assumption among our epistemic conditions. In contrast, Perea's (2012) correct beliefs

<sup>7</sup> There are some significant differences though. While Aumann and Brandenburger (1995) define knowledge as probability one belief in type-based structures, we use the standard notion of knowledge in state-based Aumann models to define common knowledge of conjectures. Also, our proofs critically build on (lexicographic) agreeing to disagree, whereas the proofs of Aumann and Brandenburger take a different route without using (standard) agreeing to disagree.

assumption essentially states that each player believes his opponent to only lexicographically deem possible the player's actual lexicographic belief hierarchy. While his epistemic operator is thus doxastic and the uncertainty is spanned by the full belief hierarchies, our correct beliefs assumption uses the stronger operator of knowledge but only concerns the players' conjectures in terms of uncertainty. Finally, Perea's (2012) notion of caution is more restrictive than ours. A player is cautious according to Perea (2012), whenever, if he lexicographically deems possible a type for any opponent, then he also lexicographically deems possible any strategy for that type. In contrast, a player already satisfies caution in our game-theoretic framework, whenever his lexicographic conjecture deems possible any strategy for all of his opponents.

## 2. Preliminaries

In state-based interactive epistemology, knowledge and beliefs are modelled within the framework of Aumann structures. Formally, an Aumann structure

$$\mathcal{A} := (\Omega, (I_i)_{i \in I}, p)$$

consists of a finite set  $\Omega$  of possible worlds (also called states of the world), a finite set  $I$  of agents, a possibility partition  $I_i$  of  $\Omega$  for every agent  $i \in I$ , and a common prior  $p : \Omega \rightarrow [0, 1]$  such that  $\sum_{\omega \in \Omega} p(\omega) = 1$ . The cell of  $I_i$  containing the world  $\omega$  is denoted by  $I_i(\omega)$  and assembles those worlds deemed possible by agent  $i$  at world  $\omega$ . It is standard to impose the so-called non-null information assumption which ensures that no information is excluded a priori, i.e.  $p(I_i(\omega)) > 0$  for all  $i \in I$  and for all  $\omega \in \Omega$ .

Agents reason about events which are defined as sets of possible worlds. The common prior  $p$  naturally extends to a measure  $p : 2^\Omega \rightarrow [0, 1]$  on the event space by setting  $p(E) = \sum_{\omega \in E} p(\omega)$  for all  $E \in 2^\Omega$ . Agents are Bayesians and consequently update the common prior with their private information as follows: the posterior belief of agent  $i$  in event  $E$  at world  $\omega$  is given by

$$p(E | I_i(\omega)) = \frac{p(E \cap I_i(\omega))}{p(I_i(\omega))}$$

and forms the decision-relevant belief of the agent.

Knowledge is formalized in terms of events. The event of agent  $i$  knowing event  $E$ , denoted by  $K_i(E)$ , is defined as

$$K_i(E) := \{\omega \in \Omega : I_i(\omega) \subseteq E\}.$$

If  $\omega \in K_i(E)$ , then  $i$  is said to know  $E$  at  $\omega$ . Mutual knowledge is given by

$$K(E) := \bigcap_{i \in I} K_i(E).$$

Setting  $K^0(E) := E$ , higher-order mutual knowledge is inductively defined by

$$K^m(E) := K(K^{m-1}(E))$$

for all  $m > 0$ . Mutual knowledge can also be denoted as 1-order mutual knowledge. The conjunction of all higher-order mutual knowledge yields common knowledge, which is formally defined as

$$CK(E) := \bigcap_{m > 0} K^m(E)$$

for all  $E \in 2^\Omega$ . This is often called the iterative definition of common knowledge. An equivalent formulation due to Aumann (1976) is based on the meet of the agents' possibility partitions and typically denoted as the meet definition of common knowledge.<sup>8</sup> Accordingly, common

knowledge is constructed as

$$CK(E) := \{\omega \in \Omega : (\bigwedge_{i \in I} I_i)(\omega) \subseteq E\}$$

for all  $E \in 2^\Omega$ , where  $(\bigwedge_{i \in I} I_i)(\omega)$  is the cell of the meet that contains the world  $\omega$ .<sup>9</sup>

Lexicographic beliefs are modelled in line with Blume et al. (1991a)'s notion of lexicographic probability systems. The following definition provides a direct adaptation of Blume et al. (1991a, Definition 3.1) to the interactive setting with multiple agents.

**Definition 1.** Let  $\Omega$  be a set of possible worlds,  $I$  be a set of agents, and  $M_i > 0$  be some integer. A *lexicographic probability system for agent  $i \in I$  ( $i$ -LPS)* is a tuple

$$\rho_i = (p_i^1, \dots, p_i^{M_i}),$$

where  $p_i^m \in \Delta(\Omega)$  for all  $m \in \{1, \dots, M_i\}$ .

Lexicographic beliefs are thus sequences of standard beliefs. The index numbers of a lexicographic probability system are also referred to as lexicographic levels.

Incorporating lexicographic probability systems into Aumann structures gives rise to the notion of lexicographic Aumann structures.

**Definition 2.** A *lexicographic Aumann structure* is a tuple

$$\mathcal{A}_L = (\Omega, I, (I_i)_{i \in I}, (\rho_i)_{i \in I}),$$

where

- $\Omega$  is a set of possible worlds,
- $I$  is a set of agents,
- $I_i \subseteq 2^\Omega$  is a possibility partition of  $\Omega$  for every agent  $i \in I$ ,
- $\rho_i = (p_i^1, \dots, p_i^{M_i})$  is an  $i$ -LPS for every agent  $i \in I$ ,
- for every agent  $i \in I$  and for every world  $\omega \in \Omega$ , there exists a lexicographic level  $m \in \{1, \dots, M_i\}$  such that  $p_i^m(I_i(\omega)) > 0$ .

The fifth item of Definition 2 ensures that no information is excluded a priori, and formally reflects the idea of caution. Actually, this condition can be seen as the lexicographic analogue to Aumann (1976)'s requirement for all information cells to be non-null events in the standard framework of Aumann structures. Caution could also be modelled as follows: for all  $i \in I$  and for all  $\omega \in \Omega$  there exists  $m \in \{1, \dots, M_i\}$  such that  $p_i^m(\omega) > 0$ . Such a condition is stronger, as it requires that every world – as opposed to only the information received – is deemed possible at some lexicographic level. The fifth item of Definition 2 is thus preferable.

Agents use their information to reason lexicographically about events. Formally, we adjust Blume et al. (1991a, Definition 4.2) to the context of lexicographic Aumann structures.

**Definition 3.** Let  $\mathcal{A}_L$  be a lexicographic Aumann structure,  $\omega \in \Omega$  be some world, and  $i \in I$  be some agent. The *conditional lexicographic probability system of agent  $i$  given his information at world  $\omega$  ( $\omega$ -conditional  $i$ -LPS)* is the tuple

$$\rho_i^\omega = (p_i^{m_1}(\cdot | I_i(\omega)), \dots, p_i^{m_L}(\cdot | I_i(\omega)))$$

where

<sup>9</sup> In fact, Brandenburger and Dekel (1987) propose a more general definition of common knowledge that can be used without the non-null information assumption holding (e.g. in situations where the set  $\Omega$  of possible worlds is uncountable). They require posterior beliefs to be proper regular conditional probabilities and modify the agents' possibility partitions appropriately in the case of null cells. Their notion of common knowledge is iterative and based on knowledge as probability one posterior belief.

<sup>8</sup> Given two partitions  $\mathcal{P}_1$  and  $\mathcal{P}_2$  of some set  $S$ , the partition  $\mathcal{P}_1$  is called *finer* than the partition  $\mathcal{P}_2$  (or  $\mathcal{P}_2$  *coarser* than  $\mathcal{P}_1$ ), if each cell of  $\mathcal{P}_1$  is a subset of some cell of  $\mathcal{P}_2$ . Given  $n$  partitions  $\mathcal{P}_1, \mathcal{P}_2, \dots, \mathcal{P}_n$  of  $S$ , the finest partition that is coarser than  $\mathcal{P}_1, \mathcal{P}_2, \dots, \mathcal{P}_n$  is called the *meet* of  $\mathcal{P}_1, \mathcal{P}_2, \dots, \mathcal{P}_n$  and is denoted by  $\bigwedge_{i=1}^n \mathcal{P}_i$ . Moreover, given  $x \in S$ , the cell of the meet  $\bigwedge_{i=1}^n \mathcal{P}_i$  containing  $x$  is denoted by  $(\bigwedge_{i=1}^n \mathcal{P}_i)(x)$ .

- the finite sequence of indices  $(m_l)_{l=0}^L$  is inductively defined by  $m_0 := 0$  and  $m_l := \min\{m \in \mathbb{N} : m_{l-1} < m \leq M_i \text{ and } p_i^m(I_i(\omega)) > 0\}$  if  $l > 0$ ;
- $p_i^{m_l}(E \mid I_i(\omega)) = \frac{p_i^{m_l}(E \cap I_i(\omega))}{p_i^{m_l}(I_i(\omega))}$  for all  $E \in 2^\Omega$  and for all  $l \in \{1, \dots, L\}$ .

An essential difference between lexicographic Aumann structures and the standard framework resides in the former equipping agents with multiple levels of – and not unique – posteriors beliefs. Technically, the sequence  $(m_l)_{l=1}^L$  of indices belonging to the  $\omega$ -conditional  $i$ -LPS  $\rho_i^\omega$  depends on both  $i$  and  $\omega$  and should thus strictly speaking be written as  $(m_{i,\omega,l})_{l=1}^{L_{i,\omega}}$ . For the sake of simplicity, the shortcut notation  $(m_l)_{l=1}^L$  is adopted, whenever the dependence on  $i$  and  $\omega$  is clear from the context. Furthermore, attention is restricted to the first  $L$  lexicographic posterior levels, where  $L := \min\{L_{i,\omega} > 0 : i \in I \text{ and } \omega \in \Omega\}$ , in order to ensure that the conditional lexicographic probability systems of every agent at every world have the same length. This restriction is only imposed for technical reasons, so that the lexicographic level posteriors the agents interactively reason about exist for all agents. Otherwise events such as “equal posteriors at all lexicographic levels” could not be properly defined. Besides, note that the lexicographic character of lexicographic probability systems actually crystallizes in two ways: an agent’s prior as well as posterior are furnished with a lexicographic structure.

The common prior assumption in Aumann structures can be directly generalized to the lexicographic setting.

**Definition 4.** Let  $\mathcal{A}_L$  be a lexicographic Aumann structure. The lexicographic Aumann structure  $\mathcal{A}_L$  satisfies the *common prior assumption (CPA)*, if there exists  $\rho = (p^1, \dots, p^M) \in (\Delta(\Omega))^M$  such that  $M = \min\{M_i \in \mathbb{N} : i \in I\}$  and  $p_i^m = p^m$  for all  $i \in I$  and for all  $m \in \{1, \dots, M\}$ . In this case, the tuple  $\rho$  is called *common prior* and  $\mathcal{A}_{LCP} = (\Omega, I, (I_i)_{i \in I}, \rho)$  is called *lexicographic Aumann structure with a common prior*.

With the existence of a common prior, the  $\omega$ -conditional  $i$ -LPS thus becomes:

$$\rho_i^\omega = \rho(\cdot \mid I_i(\omega)) = \left( p^{m_1}(\cdot \mid I_i(\omega)), \dots, p^{m_L}(\cdot \mid I_i(\omega)) \right)$$

Analogously to the case of subjective priors, the sequence  $(m_l)_{l=1}^L$  of indices should strictly speaking be written as  $(m_{i,\omega,l})_{l=1}^{L_{i,\omega}}$ , which we refrain from doing whenever the dependence on  $i$  and  $\omega$  is clear from the context.

To preempt any potential confusion about the lexicographic notation: the prior levels are denoted by  $m \in \{1, \dots, M\}$ , while the posterior levels are represented by  $l \in \{1, \dots, L\}$ . The  $l$ th posterior level corresponds to the prior level  $m_l \in \{1, \dots, M\}$  for all  $l \in \{1, \dots, L\}$ .

According to so-called Harsanyi consistency, differences in agents’ beliefs are to be attributed entirely to differences in the agents’ information. This doctrine extends to our more general set-up with lexicographic beliefs. Indeed, Definition 3 ensures that posterior heterogeneity is already excluded in the case of the common prior assumption being satisfied, if the agents face symmetric information (i.e. receive precisely the same information). Consequently, distinct posteriors need to be due to information variety.

As an illustration of our formal framework as embodied by Definitions 1 to 4, consider again the sea shipment allusion from Section 1. A lexicographic Aumann structure (cf. Definition 2) would represent a situation, where different merchants hold contingent prior beliefs and are equipped with private information about the arrival of some sea shipment. Suppose that  $\Omega = \{\omega_1, \omega_2, \omega_3, \omega_4, \omega_5, \omega_6, \omega_7, \omega_8\}$  comprises eight worlds. The eight worlds describe eight possible scenarios that are conceivable by all the merchants:

- the shipment arrives in fine weather with no pirate attack occurring ( $\omega_1 \in \Omega$ ),

- the sea shipment does not arrive in fine weather with no pirate attack occurring ( $\omega_2 \in \Omega$ ),
- the shipment arrives in adverse weather with no pirate attack occurring ( $\omega_3 \in \Omega$ ),
- the shipment does not arrive in adverse weather with no pirate attack occurring ( $\omega_4 \in \Omega$ ),
- the shipment arrives in fine weather with pirates attacking ( $\omega_5 \in \Omega$ ),
- the shipment does not arrive in fine weather with pirates attacking ( $\omega_6 \in \Omega$ ),
- the shipment arrives in adverse weather with pirates attacking ( $\omega_7 \in \Omega$ ),
- the shipment does not arrive in adverse weather with pirates attacking ( $\omega_8 \in \Omega$ ).

Suppose that some merchant  $i \in I$  deems it substantially more likely that a pirate attack does not occur. In fact, he only considers the latter to be a hypothetical contingency but he nonetheless does not discard it entirely from his thinking. Suppose further that  $i$  enjoys access to a reliable meteorological source which is signalling fine weather conditions. Such a state of mind could be modelled in our framework as follows. Merchant  $i$ ’s information partition could be given by  $I_i = \{\{\omega_1, \omega_2, \omega_5, \omega_6\}, \{\omega_3, \omega_4, \omega_7, \omega_8\}\}$  and suppose that his subjective prior would be given by an  $i$ -LPS (cf. Definition 1) as follows:  $\rho_i = (p_i^1, p_i^2)$  such that  $p_i^1(\omega_1) = \frac{4}{9}$ ,  $p_i^1(\omega_2) = p_i^1(\omega_3) = \frac{1}{9}$ , and  $p_i^1(\omega_4) = \frac{3}{9}$ , as well as  $p_i^2(\omega_5) = \frac{1}{4}$ ,  $p_i^2(\omega_6) = \frac{1}{8}$ ,  $p_i^2(\omega_7) = \frac{1}{8}$ , and  $p_i^2(\omega_8) = \frac{1}{2}$ . Assume that the shipment does arrive under fine weather conditions while withstanding a pirates’ attack. Formally speaking,  $\omega_5$  becomes the actual state of the world. The relevant posterior of merchant  $i$  is the  $\omega_5$ -conditional  $i$ -LPS (cf. Definition 3) which then obtains as  $\rho_i^{\omega_5} = p_i^{m_1}(\cdot \mid I_i(\omega_5))$ ,  $p_i^{m_2}(\cdot \mid I_i(\omega_5))$  such that  $p_i^{m_1}(\omega_1 \mid I_i(\omega_5)) = \frac{4}{5}$  and  $p_i^{m_1}(\omega_2 \mid I_i(\omega_5)) = \frac{1}{5}$ , as well as  $p_i^{m_2}(\omega_5 \mid I_i(\omega_5)) = \frac{2}{3}$  and  $p_i^{m_2}(\omega_6 \mid I_i(\omega_5)) = \frac{1}{3}$ . Moreover, in the case of the merchants being like-minded – for instance due to similar relevant past experiences with sea shipments – a common prior (cf. Definition 4) could be imposed. The sequence of prior beliefs would then be the same for all merchants, i.e. there would exist  $\rho = (p^1, \dots, p^M)$  such that  $\rho_j = \rho$  for all  $j \in I$ .

### 3. Weak agreement

Since the agents hold levels of posterior beliefs, agreement becomes a multifarious notion. Identical beliefs can obtain (or not) at different lexicographic layers. In fact, it is now shown that common knowledge of lexicographic posteriors ensures the agents’ first level posterior beliefs to coincide.

**Theorem 1 (WAT).** Let  $\mathcal{A}_{LCP}$  be a lexicographic Aumann structure with a common prior,  $E \subseteq \Omega$  be some event, and  $\omega \in \Omega$  be some world. If

$$CK\left(\bigcap_{i \in I} \bigcap_{l \in \{1, \dots, L\}} \{\omega' \in \Omega : p^{m_l}(E \mid I_i(\omega')) = p^{m_l}(E \mid I_i(\omega))\}\right) \neq \emptyset,$$

then

$$p^{m_1}(E \mid I_i(\omega)) = p^{m_1}(E \mid I_j(\omega))$$

for all  $i, j \in I$ .

**Proof.** Let  $j \in I$  be some agent,  $A_j \subseteq \Omega$  be some set such that  $(\bigwedge_{i \in I} I_i)(\omega) = \bigcup_{\omega' \in A_j} I_j(\omega')$  and  $I_j(\omega_1) \cap I_j(\omega_2) = \emptyset$  for all  $\omega_1, \omega_2 \in A_j$ . Moreover, let  $m \in \{1, \dots, M\}$  be the first lexicographic level such that  $p^m((\bigwedge_{i \in I} I_i)(\omega)) > 0$ . Consider some world  $\bar{\omega} \in A_j$ . If  $p^m(I_j(\bar{\omega})) > 0$ , then  $p^{m_1}(\cdot \mid I_j(\bar{\omega})) = p^m(\cdot \mid I_j(\bar{\omega}))$ , and by Bayesian updating,

$$p^{m_1}(E \mid I_j(\bar{\omega})) \cdot p^m(I_j(\bar{\omega})) = p^m(E \cap I_j(\bar{\omega}))$$

holds. Alternatively, if  $p^m(I_j(\bar{\omega})) = 0$ , then  $p^m(E \cap I_j(\bar{\omega})) = 0$ . Since  $p^{m_1}(\cdot \mid I_j(\bar{\omega}))$  is well-defined,

$$p^{m_1}(E \mid I_j(\bar{\omega})) \cdot p^m(I_j(\bar{\omega})) = p^m(E \cap I_j(\bar{\omega}))$$

holds trivially. Therefore,

$$p^{m_1}(E | I_j(\omega')) \cdot p^m(I_j(\omega')) = p^m(E \cap I_j(\omega'))$$

obtains for all  $\omega' \in A_j$ .

As

$$\begin{aligned} A_j &\subseteq \left( \bigcap_{i \in I} I_i(\omega) \right) \\ &\subseteq CK \left( \bigcap_{i \in I} \bigcap_{l \in \{1, \dots, L\}} \{ \omega' \in \Omega : p^{m_l}(E | I_i(\omega')) = p^{m_l}(E | I_i(\omega)) \} \right) \\ &\subseteq \bigcap_{i \in I} \bigcap_{l \in \{1, \dots, L\}} \{ \omega' \in \Omega : p^{m_l}(E | I_i(\omega')) = p^{m_l}(E | I_i(\omega)) \}, \end{aligned}$$

it is the case that  $p^{m_l}(E | I_i(\omega')) = p^{m_l}(E | I_i(\omega))$ , for all  $i \in I$  for all  $l \in \{1, \dots, L\}$  and for all  $\omega' \in A_j$ . In particular,  $p^{m_1}(E | I_j(\omega')) = p^{m_1}(E | I_j(\omega))$  holds for all  $\omega' \in A_j$ . It follows that

$$p^{m_1}(E | I_j(\omega)) \cdot p^m(I_j(\omega')) = p^m(E \cap I_j(\omega'))$$

holds for all  $\omega' \in A_j$ . Summing over all  $\omega' \in A_j$  and using countable additivity yields

$$p^{m_1}(E | I_j(\omega)) = \frac{p^m(E \cap (\bigcap_{i \in I} I_i(\omega)))}{p^m((\bigcap_{i \in I} I_i(\omega)))}.$$

Since  $j$  has been chosen arbitrarily, it can be concluded that

$$p^{m_1}(E | I_i(\omega)) = p^{m_1}(E | I_j(\omega))$$

for all  $i, j \in I$ . ■

Agents can thus not agree to disagree on their first level posterior beliefs. The preceding result remains silent though on any lexicographic level deeper than level one. In this sense, **WAT** establishes a form of weak agreement within the lexicographic framework.

Note that it is not possible to establish **WAT** by simply truncating the lexicographic Aumann structure at the first prior level and then applying Aumann's proof of his original agreement theorem to this simpler structure. This is because the first level prior may not assign positive probability to some agent's information cell, which in turn implies that a deeper level prior needs to be invoked to compute his first level posterior. Such possibilities need to be accommodated by the proof of weak agreement theorem.

For the special case of exclusively admitting the first level posteriors – formally, only considering  $p^{m_1}(\cdot | I_i(\omega))$  for all  $\omega \in \Omega$  and for all  $i \in I$  – our framework of lexicographic Aumann structures becomes essentially equivalent to [Bach and Perea \(2013\)](#)'s model, which only employs a lexicographic common prior but unique posteriors. Their non-overlapping support condition across lexicographic prior levels is not assumed in our framework though. Thus, **WAT** can be seen as a generalization of [Bach and Perea \(2013, Theorem 1\)](#). If not only the posteriors but also the common prior are restricted to a single probability measure, i.e.  $M = 1$ , then [Aumann \(1976\)](#)'s model can be recovered and **WAT** becomes the original agreement theorem.

#### 4. Disagreement

Attention is now focussed on the deeper lexicographic levels. It turns out that agents can agree to disagree on posteriors beyond the first lexicographic level.

**Proposition 1 (DIS).** *There exist a lexicographic Aumann structure  $\mathcal{A}_{LCP}$  with a common prior, some event  $E \subseteq \Omega$ , and some world  $\omega \in \Omega$ , such that*

$$CK \left( \bigcap_{i \in I} \bigcap_{l \in \{1, \dots, L\}} \{ \omega' \in \Omega : p^{m_l}(E | I_i(\omega')) = p^{m_l}(E | I_i(\omega)) \} \right) \neq \emptyset$$

and

$$p^{m_{l^*}}(E | I_i(\omega)) \neq p^{m_{l^*}}(E | I_j(\omega))$$

for some  $i, j \in I$  and for some  $l^* \in \{2, \dots, L\}$ .

**Proof.** Let  $\mathcal{A}_{LCP} = (\Omega, I, (I_i)_{i \in I}, \rho)$  be a lexicographic Aumann structure with a common prior, where

- $\Omega = \{\omega_1, \omega_2, \omega_3, \omega_4\}$ ,
- $I = \{Alice, Bob\}$ ,
- $I_{Alice} = \{\{\omega_1, \omega_2\}, \{\omega_3, \omega_4\}\}$ ,
- $I_{Bob} = \{\Omega\}$ ,
- and  $\rho = (p^1, p^2, p^3)$  with  $p^1(\omega_1) = 1$ ,  $p^2(\omega_2) = \frac{1}{3}$ ,  $p^2(\omega_3) = \frac{2}{3}$ ,  $p^3(\omega_4) = 1$ .

Consider the event  $E = \{\omega_1, \omega_3\}$ . Observe that

$$p^{m_1}(E | I_{Alice}(\omega)) = p^1(E | I_{Alice}(\omega)) = 1$$

for all  $\omega \in \{\omega_1, \omega_2\}$ , and

$$p^{m_1}(E | I_{Alice}(\omega)) = p^2(E | I_{Alice}(\omega)) = 1$$

for all  $\omega \in \{\omega_3, \omega_4\}$ .<sup>10</sup> Consequently,  $p^{m_1}(E | I_{Alice}(\omega)) = 1$  obtains at every world  $\omega \in \Omega$ . Also, observe that

$$p^{m_1}(E | I_{Bob}(\omega)) = p^1(E | I_{Bob}(\omega)) = 1$$

for all  $\omega \in \Omega$ . Therefore, *Alice's* and *Bob's* first level posterior beliefs of  $E$  coincide.

Moreover, it is the case that

$$p^{m_2}(E | I_{Alice}(\omega)) = p^2(E | I_{Alice}(\omega)) = 0$$

for all  $\omega \in \{\omega_1, \omega_2\}$ , and

$$p^{m_2}(E | I_{Alice}(\omega)) = p^3(E | I_{Alice}(\omega)) = 0$$

for all  $\omega \in \{\omega_3, \omega_4\}$ . Hence,  $p^{m_2}(E | I_{Alice}(\omega)) = 0$  obtains at every world  $\omega \in \Omega$ . Also,

$$p^{m_2}(E | I_{Bob}(\omega)) = p^2(E | I_{Bob}(\omega)) = \frac{2}{3}$$

holds at every world  $\omega \in \Omega$ . Therefore, *Alice's* and *Bob's* second level posterior beliefs of  $E$  do not coincide.

Taking  $\omega = \omega_1$  guarantees that

$$\begin{aligned} CK \left( \bigcap_{i \in I} \bigcap_{l \in \{1, \dots, L\}} \{ \omega' \in \Omega : p^{m_l}(E | I_i(\omega')) = p^{m_l}(E | I_i(\omega)) \} \right) &= CK(\Omega) \\ &= \Omega \neq \emptyset, \end{aligned}$$

while

$$p^{m_2}(E | I_{Alice}(\omega)) = 0 \neq \frac{2}{3} = p^{m_2}(E | I_{Bob}(\omega))$$

obtains at the second lexicographic level  $m_2$ . ■

A possibility result on agreeing to disagree thus emerges with lexicographic probability systems. Common knowledge of the agents' lexicographic posteriors does manifestly not suffice to establish agreement at all lexicographic levels. The agents can entertain distinct posteriors at lexicographic levels beyond one, and at the same time acknowledge this divergence. This result is somewhat surprising as it lexicographically counters Aumann's impossibility theorem. Besides, note that **DIS** would still apply and the same proof would remain valid, if a disjoint support condition were to be imposed on the lexicographic level priors.

Conceptually, **DIS** raises the question as to what drives the disagreement in a lexicographically enriched set-up. From Aumann's agreement theorem, it is typically concluded that asymmetric information does not suffice to explain heterogeneity in posterior beliefs of Bayesian agents with a common prior. Consequently, disagreement can be reached

<sup>10</sup> Recall that in the expressions  $p^{m_1}(E | I_{Alice}(\omega_1))$  and  $p^{m_1}(E | I_{Alice}(\omega_3))$ , index  $m_1$  is a shortcut notation for the two different indices  $m_{Alice, \omega_1, 1}$  and  $m_{Alice, \omega_3, 1}$ , respectively. Hence, equalities  $p^{m_1}(E | I_{Alice}(\omega_1)) = p^1(E | I_{Alice}(\omega_1))$  and  $p^{m_1}(E | I_{Alice}(\omega_3)) = p^2(E | I_{Alice}(\omega_3))$  imply that  $m_{Alice, \omega_1, 1} = 1$  and  $m_{Alice, \omega_3, 1} = 2$ , respectively.

by either weakening the common knowledge assumption or the common prior assumption. Such a conclusion does no longer apply in our lexicographic framework, since by **DIS** heterogeneous posteriors can obtain despite common knowledge of posteriors as well as the common prior remaining intact. In contrast to Aumann's original setup with standard beliefs, the lexicographic beliefs in our framework are capable of capturing hypothetical reasoning. The conceptual conclusion of Aumann's impossibility result with regard to disagreement is thus refined by **DIS** which detects hypothetical reasoning as a third source for heterogeneity in posterior beliefs.

## 5. Strong agreement

The impossibility theorem of **WAT** is weak in the sense that it only affects the first lexicographic posterior level and agreement can already fall apart at the second level as **DIS** shows. Further assumptions about the agents' like-mindedness are thus needed for a stronger result yielding equal posteriors at every lexicographic level. For this purpose an adaptation of absolute mutual absolute continuity from probability theory is introduced.

**Definition 5.** Let  $\mathcal{A}_{LCP}$  be a lexicographic Aumann structure with a common prior and  $\omega \in \Omega$  be some world. The common prior  $\rho$  is *mutually absolutely continuous*, whenever

$$p^m(I_i(\omega)) = 0, \text{ if and only if, } p^m(I_j(\omega)) = 0$$

for all  $\omega \in \Omega$ , for all  $i, j \in I$ , and for all  $m \in \{1, \dots, M\}$ .

Mutual absolute continuity ensures that at every lexicographic level the corresponding common prior handles the agents' information in synchrony. In any conceivable contingency, either the received private information at a world is deemed possible for all agents or it is excluded for everyone. Mutual absolute continuity can thus be viewed as a kind of lexicographic "same-excluding" condition.

The interpretation of the common prior assumption in the original Aumann structures with standard beliefs as agent like-mindedness can be adapted to our framework with lexicographic beliefs. The lexicographic common prior adds a contingent form of like-mindedness that also covers the different layers of hypothetical reasoning a priori. In this sense a lexicographic common prior that is mutually absolutely continuous constitutes an *intensified like-mindedness* assumption, where the players' hypothetical reasoning conditional on their information is aligned. In fact, this condition ensures that for every posterior level the agents' conditional beliefs are computed with the same level prior. If the agents violate intensified like-mindedness, then it can happen that at some posterior level they base their updated beliefs on distinct level priors. In other words, the lexicographic like-mindedness a priori gets lost in the process of Bayesian updating. The lexicographic Aumann structure constructed in the proof of **DIS** illustrates this phenomenon.

Formally, our mutual absolute continuity condition imposed on the common prior is closely related to the standard notion in probability theory which concerns two probability measures. Let  $\mu$  and  $\nu$  be measures on some set  $\Omega$ , and define  $\mu \ll \nu$ , if  $\nu(F) = 0$  implies  $\mu(F) = 0$  for all  $F \in 2^\Omega$ . Let the two measures  $\mu$  and  $\nu$  be called *standard mutually absolutely continuous*, whenever  $\mu \ll \nu$  and  $\nu \ll \mu$ .<sup>11</sup> Observe that the common prior  $\rho$  induces for every level  $m \in \{1, \dots, M\}$  and for every player  $i \in I$  a measure  $\mu_i^m : 2^\Omega \rightarrow [0, 1]$  given by

$$\mu_i^m(F) := \begin{cases} 0 & \text{if } F = \emptyset \\ \sum_{\omega \in F} \frac{p^m(I_i(\omega))}{|I_i(\omega)|} & \text{otherwise,} \end{cases}$$

<sup>11</sup> In probability theory, two mutually absolutely continuous measures are sometimes also called equivalent.

for all  $F \in 2^\Omega$ . Now, if  $\mu_i^m(F) = \sum_{\omega \in F} \frac{p^m(I_i(\omega))}{|I_i(\omega)|} > 0$  for some  $F \in 2^\Omega$ , then there exists  $\omega' \in F$  such that  $p^m(I_i(\omega')) > 0$ . By the mutual absolute continuity condition of **Definition 5**,  $p^m(I_j(\omega')) > 0$  thus holds too, and consequently  $\mu_j^m(F) = \sum_{\omega \in F} \frac{p^m(I_j(\omega))}{|I_j(\omega)|} > 0$ . Conversely, if  $p^m(I_i(\omega)) > 0$  for some  $\omega \in \Omega$ , then  $\mu_i^m(\{\omega\}) > 0$ . By standard mutual absolute continuity,  $\mu_j^m(\{\omega\}) > 0$  hence also obtains, and consequently  $p^m(I_j(\omega)) > 0$ . Therefore, the following formal characterization our mutual absolute continuity adaptation in terms of standard mutual absolute continuity from probability theory ensues.

**Remark 1.** Let  $\mathcal{A}_{LCP}$  be a lexicographic Aumann structure with a common prior. The common prior  $\rho$  is mutually absolutely continuous, if and only if,  $\mu_i^m$  and  $\mu_j^m$  are standard mutually absolutely continuous for all  $i, j \in I$  and for all  $m \in \{1, \dots, M\}$ .

Mutual absolute continuity in line with **Definition 5** can thus be viewed as a variant of standard mutual absolute continuity from probability theory.

In fact, our condition of **Definition 5** is also similar to **Stuart (1997)**'s use of mutual absolute continuity.<sup>12</sup> Accordingly, if some agent's belief assigns a positive probability to a state (which essentially corresponds to a possible world in our framework), then so do all the other agents. Even though **Stuart (1997)** does not impose any priors, an agent's belief in his model can be viewed as a posterior. While the underlying idea of **Stuart's (1997)** mutual absolute continuity and ours is the same – some form of synchronicity in both consideration and omission – his version concerns posterior beliefs and possible worlds, whereas ours refers to prior beliefs and information.

It turns out that mutual absolute continuity together with the common prior assumption and common knowledge of posteriors implies that the agents' posterior beliefs coincide at all lexicographic levels.

**Theorem 2 (SAT).** Let  $\mathcal{A}_{LCP}$  be a lexicographic Aumann structure with a common prior,  $E \subseteq \Omega$  be some event, and  $\omega \in \Omega$  be some world. If  $\rho$  is mutually absolutely continuous and

$$CK\left(\bigcap_{i \in I} \bigcap_{l \in \{1, \dots, L\}} \{\omega' \in \Omega : p^{m_l}(E | I_i(\omega')) = p^{m_l}(E | I_i(\omega))\}\right) \neq \emptyset,$$

then

$$p^{m_l}(E | I_i(\omega)) = p^{m_l}(E | I_j(\omega))$$

for all  $i, j \in I$  and for all  $l \in \{1, \dots, L\}$ .

**Proof.** We first show that if  $\rho$  is mutually absolutely continuous, then the lexicographic indices of the  $\omega'$ -conditional  $i$ -LPS  $\rho_i^{\omega'}$  are the same for all  $\omega' \in (\bigwedge_{i \in I} I_i(\omega))$  and for all  $i \in I$ . Let  $j \in I$ ,  $\omega' \in (\bigwedge_{i \in I} I_i(\omega))$  as well as  $(m_l)_{l=1}^L$  and  $(m'_l)_{l=1}^L$  be the indices of  $\rho_j^{\omega'}$  and  $\rho_j^{\omega'}$ , respectively. Since  $\omega' \in (\bigwedge_{i \in I} I_i(\omega))$ , the world  $\omega'$  is doxastically reachable from  $\omega$ , i.e., there exists a sequence  $(P^k)_{k=1}^N$  of information cells such that  $\omega \in P^1$ ,  $\omega' \in P^N$ , and  $P^k \cap P^{k+1} \neq \emptyset$  for all  $1 \leq k < N$ . Since  $\rho$  is mutually absolutely continuous, it is the case that,  $p^m(P^k) = 0$  if and only if  $p^m(P^{k+1}) = 0$  for all  $m \in \{1, \dots, M\}$  and for all  $1 \leq k < N$ . Thus,  $p^m(P^1) = 0$  if and only if  $p^m(P^N) = 0$  for all  $m \in \{1, \dots, M\}$ . Since  $\omega \in I_j(\omega) \cap P^1$ ,  $\omega' \in I_j(\omega') \cap P^N$  and  $\rho$  is mutually absolutely continuous, it follows that  $p^m(I_j(\omega)) = 0$  if and only if  $p^m(P^1) = 0$  and  $p^m(I_j(\omega')) = 0$  if and only if  $p^m(P^N) = 0$ , and thus  $p^m(I_j(\omega)) = 0$  if and only if  $p^m(I_j(\omega')) = 0$ , for all  $m \in \{1, \dots, M\}$ . Consequently,  $(m_l)_{l=1}^L = (m'_l)_{l=1}^L$ . Now, towards a contradiction, suppose that there exist  $j' \in I$  and  $l \in \{1, \dots, L\}$  such that  $m'_l \neq m''_l$ , where  $(m''_l)_{l=1}^L$  are the indices of  $\rho_{j'}^{\omega'}$ . Without loss of generality, suppose that  $l$  is the least such index. Then,

<sup>12</sup> In **Stuart (1997)**, mutual absolute continuity plays an important role in establishing all period defection in the normal-form model of the finitely repeated prisoners' dilemma.

either  $m'_l < m''_l$ , in which case,  $p^{m'_l}(I_j(\omega')) > 0$  and  $p^{m''_l}(I_{j'}(\omega')) = 0$ , or  $m'_l > m''_l$ , in which case,  $p^{m'_l}(I_j(\omega')) = 0$  and  $p^{m''_l}(I_{j'}(\omega')) > 0$ . In both cases, a contradiction with the mutual absolute continuity of  $\rho$  obtains. Consequently,  $(m_l)_{l=1}^L = (m'_l)_{l=1}^L = (m''_l)_{l=1}^L =: (\bar{m}_l)_{l=1}^L$ . The  $\omega'$ -conditional  $i$ -LPS can then be written as

$$\rho_i^{\omega'} = \rho(\cdot | I_i(\omega')) = (p^{\bar{m}_1}(\cdot | I_i(\omega')), \dots, p^{\bar{m}_L}(\cdot | I_i(\omega')))$$

for all  $i \in I$  and for all  $\omega' \in (\bigwedge_{i \in I} I_i)(\omega)$ .

We are now ready to derive agreement in posteriors. Let  $j' \in I$  and  $A_{j'} \subseteq \Omega$  such that  $(\bigwedge_{i \in I} I_i)(\omega) = \bigcup_{\omega' \in A_{j'}} I_{j'}(\omega')$  and  $I_{j'}(\omega_1) \cap I_{j'}(\omega_2) = \emptyset$  for all  $\omega_1, \omega_2 \in A_{j'}$ . Note that

$$\begin{aligned} A_{j'} &\subseteq \left( \bigwedge_{i \in I} I_i \right)(\omega) \\ &\subseteq CK \left( \bigcap_{i \in I} \bigcap_{l \in \{1, \dots, L\}} \{ \omega' \in \Omega : p^{m_l}(E | I_i(\omega')) = p^{m_l}(E | I_i(\omega)) \} \right) \\ &\subseteq \bigcap_{i \in I} \bigcap_{l \in \{1, \dots, L\}} \{ \omega' \in \Omega : p^{m_l}(E | I_i(\omega')) = p^{m_l}(E | I_i(\omega)) \}. \end{aligned}$$

Consider some  $l' \in \{1, \dots, L\}$ . It follows that

$$\begin{aligned} p^{m_{l'}}(E | I_{j'}(\omega)) &= p^{m_{l'}}(E | I_{j'}(\omega')) \\ &= \frac{p^{m_{l'}}(E \cap I_{j'}(\omega'))}{p^{m_{l'}}(I_{j'}(\omega'))} = \frac{p^{m_{l'}}(E \cap I_{j'}(\omega'))}{p^{m_{l'}}(I_{j'}(\omega'))} \end{aligned}$$

for all  $\omega' \in A_{j'}$ . Consequently,

$$p^{m_{l'}}(E | I_{j'}(\omega)) \cdot p^{\bar{m}_{l'}}(I_{j'}(\omega')) = p^{\bar{m}_{l'}}(E \cap I_{j'}(\omega')),$$

for all  $\omega' \in A_{j'}$ . Summing over all  $\omega' \in A_{j'}$  and using countable additivity yields

$$p^{m_{l'}}(E | I_{j'}(\omega)) = \frac{p^{\bar{m}_{l'}}(E \cap (\bigwedge_{i \in I} I_i)(\omega))}{p^{\bar{m}_{l'}}((\bigwedge_{i \in I} I_i)(\omega))} = p^{\bar{m}_{l'}}(E | (\bigwedge_{i \in I} I_i)(\omega)).$$

Since  $j'$  and  $l'$  have been chosen arbitrarily, it can be concluded that

$$p^{m_l}(E | I_i(\omega)) = p^{m_l}(E | I_j(\omega))$$

for all  $i, j \in I$  and for all  $l \in \{1, \dots, L\}$ . ■

It is thus impossible for lexicographically-minded agents to agree to disagree whenever mutual absolute continuity is satisfied. In contrast to **WAT**, which only ensures a weak form of agreement at the first posterior level, **SAT** establishes strong agreement at all lexicographic posterior levels.

From a conceptual perspective, agreement is only ensured in the lexicographically enriched framework by a substantial strengthening of the agents' like-mindedness. It does not suffice to require a common prior at all reasoning levels. On top of that, each of these priors also has to synchronically consider or synchronically neglect the agents' information in order to reconcile their updating. Together with common knowledge of posteriors, the assumption of intensified like-mindedness drives the homogeneity of the posteriors in our lexicographic framework.

The particular lexicographic Aumann structure constructed in the proof of **DIS** suggests that **SAT** qualifies as tight with respect to the mutual absolute continuity condition.<sup>13</sup> There the other two key assumptions, i.e. common prior as well common knowledge of posteriors, but not mutual absolute continuity hold, while the consequent, i.e. lexicographically identical posterior beliefs, fails.

Continuity in agreeing to lexicographically disagree follows from **SAT** in the sense that equal prior beliefs up to some lexicographic prior level imply equal posterior beliefs up to a corresponding lexicographic posterior level. Suppose that the common prior assumption is weakened

such that the agents' priors coincide up to some level  $\bar{M} < M$ , and modify the initial lexicographic Aumann structure by truncating the agent's lexicographic priors at  $\bar{M}$ , which is equivalent to imposing a common prior  $\rho = (p^1, \dots, p^{\bar{M}})$ . By **SAT** it follows that common knowledge of lexicographic posteriors at some world  $\omega \in \Omega$  implies equal posterior measures for every level  $l \in \{1, \dots, \min\{L_{i,\omega} \in \mathbb{N} : i \in I\}\}$  in the truncated structure, and hence also up to level  $\min\{L_{i,\omega} \in \mathbb{N} : i \in I\}$  in the initial lexicographic Aumann structure. In this sense, the lexicographic impossibility result of **SAT** is continuous.

## 6. Perfect equilibrium

Next, we turn to game theory where some of our results on lexicographic agreeing to disagree are employed for an epistemic analysis of tremble equilibria. In game theory, strategic interaction of multiple agents is modelled, and possible outcomes are predicted based on different assumptions. Static games with complete information constitute the most elementary analytical framework. Formally, such games are represented by a tuple

$$\Gamma = (I, (S_i)_{i \in I}, (U_i)_{i \in I})$$

consisting of a finite set  $I$  of players and finite non-empty strategy sets  $S_i$  as well as real-valued utility functions  $U_i$  with domain  $\prod_{j \in I} S_j$  for every player  $i \in I$ . In terms of notation, the set  $S_{-i} := \prod_{j \in I \setminus \{i\}} S_j$  refers to the product set of the  $i$ 's opponents' strategy combinations. The tuple  $\Gamma = (I, (S_i)_{i \in I}, (U_i)_{i \in I})$  is often also referred to as normal form. As background hypotheses it is stipulated that all players choose their strategies simultaneously and that the ingredients of the game, i.e. the normal form, is common knowledge among the players. Solution concepts propose plausibility criteria or decision rules in line with which the players are supposed to act. Formally, a solution concept defines a subset  $SC \subseteq \prod_{i \in I} S_i$  of the set of all strategy combinations as possible outcomes of the game.

The solution concept of Nash equilibrium – due to [Nash \(1950, 1951\)](#) – requires players to choose utility maximizing against fixed strategies of the opponents. In order to ensure existence of an equilibrium point in any game, also randomizations over strategies are admitted. The set of choice objects for every player  $i \in I$  is thus enlarged from  $S_i$  to  $\Delta(S_i)$ , where a typical element  $\sigma_i$  of  $\Delta(S_i)$  is called a mixed strategy of player  $i$ . The utility functions  $U_i$  are extended from  $\prod_{j \in I} S_j$  to  $\prod_{j \in I} (\Delta(S_j))$  for every player  $i \in I$  by an expected utility computation. A tuple of mixed strategies  $\sigma = (\sigma_j)_{j \in I}$  constitutes a *Nash equilibrium*, whenever

$$s_i^* \in \arg \max_{s_i \in S_i} \left\{ \sum_{s_{-i} \in S_{-i}} \left( \bigotimes_{j \in I \setminus \{i\}} \sigma_j \right) (s_{-i}) \cdot U_i(s_i, s_{-i}) \right\} \quad (1)$$

for all  $s_i^* \in \text{supp}(\sigma_i)$  and for all  $i \in I$ .<sup>14</sup> If Eq. (1) holds,  $s_i^*$  is called a best response to  $\sigma_{-i}$ , where  $\sigma_{-i} := (\sigma_j)_{j \in I \setminus \{i\}}$ . Player  $i$  is said to strictly prefer a strategy  $s_i$  to some other strategy  $s'_i$  given  $\sigma_{-i}$ , whenever  $\sum_{s_{-i} \in S_{-i}} (\bigotimes_{j \in I \setminus \{i\}} \sigma_j)(s_{-i}) \cdot U_i(s_i, s_{-i}) > \sum_{s_{-i} \in S_{-i}} (\bigotimes_{j \in I \setminus \{i\}} \sigma_j)(s_{-i}) \cdot U_i(s'_i, s_{-i})$  holds.

In classical game theory, the multiplicity of Nash equilibria in many games has been deemed unsatisfactory and refinements have thus been sought. A particular class of equilibrium refinements is based on the idea that players can make mistakes with small probability. Phrased in more vivid terms: players possibly tremble when implementing their strategies. In line with this intuition, various tremble equilibria have been proposed in the literature. The most basic such solution concept is

<sup>13</sup> Tightness is interpreted in the style of [Aumann and Brandenburger \(1995\)](#), i.e. whether dropping only one assumption of a result were to already break its conclusion.

<sup>14</sup> Given a probability measure  $p \in \Delta(X)$  on some set  $X$  its support is defined as  $\text{supp}(p) := \{x \in X : p(x) > 0\}$ . Fixing  $K \in \mathbb{N}$  and probability measures  $p_k$  on sets  $X_k$  for all  $k \in \{1, \dots, K\}$ ,  $\bigotimes_{k \in \{1, \dots, K\}} p_k$  denotes the product measure on the set  $\prod_{k \in \{1, \dots, K\}} X_k$ .

	y	z
a	1, 1	0, 0
b	0, 0	0, 0

Fig. 2. Another two player game.

Selten's (1975) perfect equilibrium.<sup>15</sup> Essentially, attention is restricted to Nash equilibria that obtain as limits of sequences of perturbed strategy combinations. While originally introduced by Selten (1975, Section 8) as a solution concept for dynamic games, perfect equilibrium has also been widely used in static games. A formal definition of perfect equilibrium for the class of static games ensues as follows.

**Definition 6.** Let  $\Gamma$  be a game and  $\sigma = (\sigma_i)_{i \in I} \in \prod_{i \in I} \Delta(S_i)$  be a tuple of mixed strategies. The tuple  $\sigma$  constitutes a *perfect equilibrium* of  $\Gamma$ , if there exists a sequence of tuples of mixed strategies  $(\sigma^k)_{k \in \mathbb{N}} = ((\sigma_i^k)_{i \in I})_{k \in \mathbb{N}} \in (\prod_{i \in I} \Delta(S_i))^{\mathbb{N}}$  such that

- (i)  $\lim_{k \rightarrow \infty} \sigma^k = \sigma$ ;
- (ii) for all  $i \in I$  and for all  $k \in \mathbb{N}$ , it is the case that  $\text{supp}(\sigma_i^k) = S_i$ ;
- (iii) for all  $i \in I$  and for all  $k \in \mathbb{N}$ , if  $s_i \in \text{supp}(\sigma_i)$ , then  $s_i$  is a best response to  $\sigma_{-i}^k$ .

A perfect equilibrium thus always coincides with the limit of a sequence of trembles. Moreover, for every player, his perfect equilibrium mixed strategy only assigns positive probability to strategies that are best responses to any of the opponents' tremble combinations. It can be shown that a perfect equilibrium must be a Nash equilibrium (Selten, 1975, Lemma 9). This result essentially rests on the fact that the expected utilities are continuous in mixed strategy profiles. Conversely, Nash equilibrium does not imply perfect equilibrium. The latter solution concept thus is stronger than the former. In classical parlance, perfect equilibrium constitutes a refinement of Nash equilibrium.

The following example illustrates these two solution concepts.

**Example 1.** Consider the two player game depicted in Fig. 2 with players Alice and Bob, where Alice chooses a "row" ( $a$  or  $b$ ) and Bob picks a "column" ( $y$  or  $z$ ). The mixed strategy tuple  $\sigma = (\sigma_{\text{Alice}}, \sigma_{\text{Bob}})$ , where  $\sigma_{\text{Alice}}(a) = 1$  and  $\sigma_{\text{Bob}}(y) = 1$ , forms a Nash equilibrium, as  $a$  is a best response to  $\sigma_{\text{Bob}}$  and  $y$  is a best response to  $\sigma_{\text{Alice}}$ . To see that  $\sigma$  also constitutes a perfect equilibrium, construct a sequence of tuples of mixed strategies  $(\sigma^k)_{k \in \mathbb{N} \setminus \{0\}} = ((\sigma_{\text{Alice}}^k, \sigma_{\text{Bob}}^k))_{k \in \mathbb{N} \setminus \{0\}}$  by setting  $\sigma_{\text{Alice}}^k(a) = 1 - \frac{1}{k+1}$ ,  $\sigma_{\text{Alice}}^k(b) = 0 + \frac{1}{k+1}$ ,  $\sigma_{\text{Bob}}^k(y) = 1 - \frac{1}{k+1}$  and  $\sigma_{\text{Bob}}^k(z) = 0 + \frac{1}{k+1}$  for all  $k \in \mathbb{N} \setminus \{0\}$ . Observe that  $\lim_{k \rightarrow \infty} \sigma^k = \sigma$  as well as  $\text{supp}(\sigma_{\text{Alice}}^k) = S_{\text{Alice}}$  and  $\text{supp}(\sigma_{\text{Bob}}^k) = S_{\text{Bob}}$  for all  $k \in \mathbb{N} \setminus \{0\}$ . Moreover,  $a$  is a best response to  $\sigma_{\text{Bob}}^k$  for all  $k \in \mathbb{N} \setminus \{0\}$  and  $y$  is a best response to  $\sigma_{\text{Alice}}^k$  for all  $k \in \mathbb{N} \setminus \{0\}$ . It follows that  $\sigma$  is a perfect equilibrium.

The mixed strategy tuple  $\sigma' = (\sigma'_{\text{Alice}}, \sigma'_{\text{Bob}})$ , where  $\sigma'_{\text{Alice}}(b) = 1$  and  $\sigma'_{\text{Bob}}(z) = 1$  also constitutes a Nash equilibrium, since  $b$  is a best response to  $\sigma_{\text{Bob}}$  and  $z$  is a best response to  $\sigma_{\text{Alice}}$ . However, it does not form a perfect equilibrium. Suppose that there exists a sequence of full support mixed strategy tuples  $(\sigma^k_{\text{Alice}}, \sigma^k_{\text{Bob}})_{k \in \mathbb{N} \setminus \{0\}} \in (\Delta(S_{\text{Alice}}) \times \Delta(S_{\text{Bob}}))^{\mathbb{N} \setminus \{0\}}$  with limit point  $\sigma'$ . Then,  $b$  cannot be a best response to  $\sigma^k_{\text{Bob}}$  for any  $k \in \mathbb{N} \setminus \{0\}$ . Indeed, as soon as  $y$  receives positive probability, only  $a$  can be a best response for Alice. It follows that  $\sigma'$  is not a perfect equilibrium. ♣

## 7. Lexicographic characterization

It is known that tremble equilibria with their sequences of full support mixed strategy tuples can be characterized in terms of lexicographic conjectures. The latter can be modelled as lexicographic probability systems in which for every player the set of opponents' choice combinations defines the basic space of uncertainty. Perfect equilibrium and proper equilibrium have been reformulated with lexicographic conjectures by Blume et al. (1991b) and shown to be equivalent to their notion of lexicographic Nash equilibrium plus further restrictions, respectively. In this section we define lexicographic perfect equilibrium and lexicographic semi-perfect equilibrium. While these two solution concepts phrased in terms of lexicographic conjectures essentially correspond to variants of Blume et al.'s (1991b) lexicographic Nash equilibrium, our definitions are aligned with our formal framework and formulated in a direct way.

Some further concepts and notation need to be introduced. Let  $\Gamma$  be a game and  $i \in I$  be some player. A sequence  $\beta_i = (b_i^1, \dots, b_i^L) \in (\Delta(S_{-i}))^L$  of probability measures, for some  $L \in \mathbb{N}$ , is called player  $i$ 's *lexicographic conjecture*. For the sake of simplicity we assume the same number  $L$  of levels for all  $i \in I$ . A lexicographic conjecture  $\beta_i$  is *cautious*, whenever for all  $j \in I \setminus \{i\}$  and for all  $s_j \in S_j$ , there exists some lexicographic level  $l^* \in \{1, \dots, L\}$  such that  $\text{marg}_{S_j} b_i^{l^*}(s_j) > 0$ , where  $\text{marg}_{S_j} b_i^{l^*}(s_j) := \sum_{s_{-i,j} \in S_{-i,j}} b_i^{l^*}(s_{-i,j}, s_j)$  for all  $s_j \in S_j$ . Given a strategy  $s_i \in S_i$  and a lexicographic conjecture  $\beta_i = (b_i^1, \dots, b_i^L) \in (\Delta(S_{-i}))^L$ ,

$$u_i^l(s_i, \beta_i) := \sum_{s_{-i} \in S_{-i}} b_i^l(s_{-i}) \cdot U_i(s_i, s_{-i})$$

is player  $i$ 's *level- $l$  expected utility* for all  $l \in \{1, \dots, L\}$ . Equipped with a lexicographic conjecture  $\beta_i \in (\Delta(S_{-i}))^L$ , player  $i$  *strictly lex-prefers* a strategy  $s_i \in S_i$  to some other strategy  $s'_i \in S_i$ , whenever there exists a lexicographic level  $l^* \in \{1, \dots, L\}$  such that

$$u_i^{l^*}(s_i, \beta_i) > u_i^{l^*}(s'_i, \beta_i) \text{ and } u_i^l(s_i, \beta_i) = u_i^l(s'_i, \beta_i)$$

for all  $l < l^*$ . A strategy  $s_i^* \in S_i$  is called *lex-optimal* given  $\beta_i$ , if there exists no strategy  $s_i \in S_i$  such that  $i$  strictly lex-prefers  $s_i$  to  $s_i^*$ . Similarly, player  $i$  is said to be *lex-indifferent* between  $s_i$  and  $s'_i$ , whenever  $u_i^l(s_i, \beta_i) = u_i^l(s'_i, \beta_i)$  for all  $l \in \{1, \dots, L\}$ . Player  $i$  *weakly lex-prefers*  $s_i$  to  $s'_i$ , if he strictly lex-prefers the former to the latter or feels lex-indifferent. A lexicographic conjecture  $\beta_i$  is called *lexicographic product conjecture*, if  $b_i^l = \bigotimes_{j \in I \setminus \{i\}} \text{marg}_{S_j} b_i^l$  holds for all  $l \in \{1, \dots, L\}$ , and is formally written as

$$\beta_i := \bigotimes_{j \in I \setminus \{i\}} \text{marg}_{S_j} \beta_i := \left( \bigotimes_{j \in I \setminus \{i\}} \text{marg}_{S_j} b_i^1, \dots, \bigotimes_{j \in I \setminus \{i\}} \text{marg}_{S_j} b_i^L \right).$$

Conceptually, a player with a lexicographic product conjecture treats his opponents' choices as uncorrelated.<sup>16</sup>

Selten's (1975) solution concept of perfect equilibrium can be expressed in terms of lexicographic conjectures.

**Definition 7.** Let  $\Gamma$  be a finite game,  $\sigma = (\sigma_i)_{i \in I} \in \prod_{i \in I} (\Delta(S_i))$  be a tuple of mixed strategies, and  $L \in \mathbb{N}$ . The tuple  $\sigma$  constitutes a *lexicographic perfect equilibrium* of  $\Gamma$ , if there exist a tuple  $\beta = (\beta_i)_{i \in I} \in \left( (\Delta(S_{-i}))^L \right)_{i \in I}$  of lexicographic conjectures and a lexicographic product measure  $\pi = (\pi^1, \dots, \pi^L) \in (\Delta(\prod_{i \in I} S_i))^L$  such that for all  $i \in I$ , the following properties hold:

- (a)  $\beta_i = (b_i^1, \dots, b_i^L)$  is cautious;

<sup>15</sup> Other tremble equilibria are, for instance, Myerson's (1978) proper equilibrium, van Damme's (1984) quasi-perfect equilibrium, as well as Harsanyi and Selten's (1988) uniformly perfect equilibrium.

<sup>16</sup> While players by assumption do choose independently of course, it is well known that this does not preclude the possibility that beliefs about opponents' choices violate statistical independence. Essentially, the reason lies in two distinct forms of independence – causal and epistemic – which do not imply each other.

- (b)  $\sigma_i = \arg_{S_i} b_j^1$  for all  $j \in I \setminus \{i\}$ ;
- (c) if  $s_i \in \text{supp}(\sigma_i)$ , then  $s_i$  is lex-optimal given  $\beta_i$ ;
- (d)  $\beta_i = \bigotimes_{j \in I \setminus \{i\}} \arg_{S_j} \beta_j$ ;
- (e)  $\arg_{S_{-i}} \pi = \beta_i$ .

A lexicographic formulation of perfect equilibrium thus builds on an interpretation of mixed strategies as conjectures. In this regard, condition (b) blocks any doxastic ambiguity by requiring that for a given player all opponents share the same belief about his choice. The trembles of the classical definition are mimicked via condition (a) which requires the lexicographic conjectures to be cautious. The best response property of the perfect equilibrium tuple is ensured by condition (c) according to which only choices supported by the player's lexicographic conjecture receive positive probability. Epistemic independence is built in via condition (d) postulating that the players' lexicographic conjectures are the product of their marginals. Each of the lexicographic beliefs are required by condition (e) to stem from a joint source. In essence, lexicographic perfect equilibrium corresponds to Blume et al.'s (1991a) lexicographic Nash equilibrium plus full support at all lexicographic levels, a common prior, and some independence condition.

The classical and the lexicographic versions of perfect equilibrium are equivalent.

**Lemma 1.** Let  $\Gamma$  be a finite game and  $\sigma \in \times_{i \in I} (\Delta(S_i))$  be a tuple of mixed strategies. The tuple  $\sigma$  constitutes a perfect equilibrium of  $\Gamma$ , if and only if,  $\sigma$  constitutes a lexicographic perfect equilibrium of  $\Gamma$ .

**Proof.** See Appendix.

The classical formulation (Section 6) and the lexicographic variant (Definition 7) of perfect equilibrium can thus be used interchangeably. Lemma 1 is by and large equivalent to Blume et al. (1991b, Proposition 7), where classical perfect equilibrium is characterized in terms of their notion of lexicographic Nash equilibrium plus some additional assumptions. For the sake of completeness and self-containedness we explicitly show the equivalence. However, since Lemma 1 lies outside the focus of this paper its proof is deferred to the Appendix.

A possibly meaningful weakening of lexicographic perfect equilibrium would obtain, if conditions (d) and (e) of Definition 7 were to be dropped.

**Definition 8.** Let  $\Gamma$  be a finite game,  $\sigma = (\sigma_i)_{i \in I} \in \times_{i \in I} (\Delta(S_i))$  be a tuple of mixed strategies, and  $L \in \mathbb{N}$ . The tuple  $\sigma$  constitutes a *lexicographic semi-perfect equilibrium* of  $\Gamma$ , if there exists a tuple  $\beta = (\beta_i)_{i \in I} \in \left( \left( \Delta(S_{-i}) \right)^L \right)_{i \in I}$  of lexicographic conjectures such that for all  $i \in I$ , the following properties hold:

- (a)  $\beta_i = (b_i^1, \dots, b_i^L)$  is cautious;
- (b)  $\sigma_i = \arg_{S_i} b_j^1$  for all  $j \in I \setminus \{i\}$ ;
- (c) if  $s_i \in \text{supp}(\sigma_i)$ , then  $s_i$  is lex-optimal given  $\beta_i$ .

A lexicographic semi-perfect equilibrium does admit a player's lexicographic conjecture about his opponents' choices to not be independent. Accordingly, he may deem it lexicographically possible for some opponents' choices to be correlated. Note that correlated beliefs at some level do not imply the belief that players do not choose independently from each other. Even though the actions of any two players in a static game are entirely autonomous, the reasoning leading to these actions might be related in a way that makes them correlated from the perspective of a third player. Also, in contrast to perfect equilibrium, more flexibility about the lexicographic conjectures is permitted by Definition 8, as they no longer need to be projections of a joint source. The solution concept of lexicographic semi-perfect equilibrium basically coincides with Blume et al.'s (1991a) notion of lexicographic Nash equilibrium plus some full support property.

It is clear that perfect equilibrium implies semi-perfect equilibrium, as the latter requires two properties less than the former. The following example shows that the converse does not hold though.

**Example 2.** Consider the three player game depicted in Fig. 3 with players Alice, Bob, and Claire, where Alice chooses a "row" ( $a$  or  $b$ ), Bob picks a "column" ( $y$  or  $z$ ), and Claire selects a "matrix" (*left*, *middle*, or *right*).

It is first shown that the mixed strategy tuple  $\sigma = (\sigma_{\text{Alice}}, \sigma_{\text{Bob}}, \sigma_{\text{Claire}})$ , where  $\sigma_{\text{Alice}}(a) = \sigma_{\text{Alice}}(b) = 0.5$ ,  $\sigma_{\text{Bob}}(y) = \sigma_{\text{Bob}}(z) = 0.5$ , and  $\sigma_{\text{Claire}}(\text{middle}) = 1$  forms a lexicographic semi-perfect equilibrium. Define conjectures  $\beta_{\text{Alice}} = (b_{\text{Alice}}^1, b_{\text{Alice}}^2)$ ,  $\beta_{\text{Bob}} = (b_{\text{Bob}}^1, b_{\text{Bob}}^2)$ , and  $\beta_{\text{Claire}} = (b_{\text{Claire}}^1, b_{\text{Claire}}^2)$  such that

$$b_{\text{Alice}}^1 = 0.5 \cdot (y, \text{middle}) + 0.5 \cdot (z, \text{middle}),$$

$$b_{\text{Alice}}^2 = 0.5 \cdot (y, \text{left}) + 0.5 \cdot (z, \text{right}),$$

$$b_{\text{Bob}}^1 = 0.5 \cdot (a, \text{middle}) + 0.5 \cdot (b, \text{middle}),$$

$$b_{\text{Bob}}^2 = 0.5 \cdot (a, \text{left}) + 0.5 \cdot (b, \text{right}),$$

$$b_{\text{Claire}}^1 = 0.5 \cdot (a, y) + 0.5 \cdot (b, z),$$

$$b_{\text{Claire}}^2 = 1 \cdot (a, y).$$

Each of the three conjectures is cautious, as all choices of all respective opponents' receive positive probability at some lexicographic level. Moreover,

$$\arg_{S_{\text{Alice}}} b_{\text{Bob}}^1 = \arg_{S_{\text{Alice}}} b_{\text{Claire}}^1 = 0.5 \cdot a + 0.5 \cdot b = \sigma_{\text{Alice}},$$

$$\arg_{S_{\text{Bob}}} b_{\text{Alice}}^1 = \arg_{S_{\text{Bob}}} b_{\text{Claire}}^1 = 0.5 \cdot y + 0.5 \cdot z = \sigma_{\text{Bob}},$$

$$\arg_{S_{\text{Claire}}} b_{\text{Alice}}^1 = \arg_{S_{\text{Claire}}} b_{\text{Bob}}^1 = \text{middle} = \sigma_{\text{Claire}}.$$

Observe that  $a$  and  $b$  are lex-optimal given  $\beta_{\text{Alice}}$ ,  $y$  and  $z$  are lex-optimal given  $\beta_{\text{Bob}}$ , as well as *middle* is lex-optimal given  $\beta_{\text{Claire}}$ . Consequently,  $\sigma$  constitutes a lexicographic semi-perfect equilibrium. However,  $\sigma$  is not lexicographic perfect, as  $b_{\text{Claire}}^1$ 's probability measure violates independence and property (d) of Definition 7 is thus not satisfied. ♣

It could be interesting to explore new solution concepts based on various weakenings of lexicographic perfect equilibrium such as lexicographic semi-perfect equilibrium. Another possibility would be to also admit conjectures that violate the projection property on the first lexicographic level. A corresponding perfect equilibrium variant could then be defined directly in terms of lexicographic conjectures and be required to satisfy the conditions (a) and (c) of Definition 7. We leave such thoughts for further research.

## 8. Epistemic characterization

We now explore the interactive reasoning assumptions of perfect equilibrium and thereby extend the work of Blume et al. (1991b). While Blume et al. (1991b) characterize perfect equilibrium in terms of lexicographic conjectures, they do not perform any epistemic analysis involving higher-order beliefs to unveil the interactive thinking that drive players to choose in line with this solution concept. The latter is precisely the focus of this section. A key role will be played by our results on lexicographic agreeing to disagree. In particular, the weak agreement theorem (WAT) as well as the strong agreement theorem (SAT) turn into essential ingredients to establish an epistemic foundation for perfect equilibrium.

In game theory, reasoning is captured by means of epistemic structures that are added to the formal framework. Different patterns or assumptions about reasoning can then be expressed in the form of epistemic hypotheses. Classical solution concepts can be characterized in terms of reasoning by epistemic conditions. In this way, the interactive thinking a solution concept requires on behalf of the players so that they act in line with its prediction is made explicit.

	y	z		y	z		y	z
a	1, 1, 2	0, 0, 0	a	1, 1, 2	0, 0, 0	a	1, 1, 0	0, 0, 0
b	0, 0, 0	1, 1, 0	b	0, 0, 0	1, 1, 2	b	0, 0, 0	1, 1, 2
	left			middle			right	

Fig. 3. A three player game.

Before we turn to reasoning foundations, some more formal structure and notions have to be fixed. First of all, the basic framework of games as embodied by  $\Gamma$  needs to be enlarged by an epistemic dimension. To this end we introduce the notion of a lexicographic Aumann model.

**Definition 9.** Let  $\Gamma$  be a finite game. A *lexicographic Aumann model* of  $\Gamma$  is a tuple

$$\mathcal{A}_{LCP}^{\Gamma} = (\Omega, \rho, I, (I_i, \hat{s}_i)_{i \in I})$$

where

- $\Omega$  is a set of possible worlds,
- $\rho = (p^1, \dots, p^M)$  is a common prior,
- $I$  is the set of players from  $\Gamma$ ,
- $I_i \subseteq 2^{\Omega}$  is player  $i$ 's possibility partition of  $\Omega$  for all  $i \in I$ ,
- $\hat{s}_i : \Omega \rightarrow S_i$  is player  $i$ 's choice function that is  $I_i$ -measurable for all  $i \in I$ , i.e.,  $\hat{s}_i(w') = \hat{s}_i(w)$  for all  $w, w' \in \Omega$  such that  $w' \in I_i(w)$ ,
- for every player  $i \in I$  and for every world  $\omega \in \Omega$ , there exists a level  $m \in \{1, \dots, M\}$  such that  $p^m(I_i(\omega)) > 0$ .

A lexicographic Aumann models thus corresponds to a lexicographic Aumann structure (Definition 2) supplemented by choice functions for every player that connect the interactive epistemology to games. It then becomes possible to express game-theoretic events and interactive beliefs as well as knowledge about these.

The event that player  $i$  chooses strategy  $s_i \in S_i$  is formalized as

$$[s_i] := \{\omega \in \Omega : \hat{s}_i(\omega) = s_i\}$$

and the event that  $i$ 's opponents choose  $s_{-i} \in S_{-i}$  is given by

$$[s_{-i}] := \bigcap_{j \in I \setminus \{i\}} [s_j].$$

Note that the  $I_i$ -measurability of  $\hat{s}_i$  ensures that either  $I_i(\omega) \subseteq [s_i]$  or  $I_i(\omega) \subseteq [s_i]^c$ . A lexicographic conjecture function can be defined as  $\hat{\beta}_i : \Omega \rightarrow (\Delta(S_{-i}))^L$ , where

$$\begin{aligned} \hat{\beta}_i(\omega)(s_{-i}) &= (\hat{\beta}_i^1(\omega)(s_{-i}), \dots, \hat{\beta}_i^L(\omega)(s_{-i})) \\ &:= \rho([s_{-i}] \mid I_i(\omega)) = (p^{m_1}([s_{-i}] \mid I_i(\omega)), \dots, p^{m_L}([s_{-i}] \mid I_i(\omega))) \end{aligned}$$

for all  $\omega \in \Omega$  and for all  $s_{-i} \in S_{-i}$ . From the  $I_i$ -measurability of the level posteriors it follows that  $\hat{\beta}_i$  is  $I_i$ -measurable too, i.e.  $\hat{\beta}_i(\omega') = \hat{\beta}_i(\omega)$  for all  $\omega' \in I_i(\omega)$ . Hence, for every lexicographic conjecture  $\beta_i$  of player  $i$ , the lexicographic conjecture function induces a coarsening of  $I_i$  of the form

$$[\beta_i] := \{\omega \in \Omega : \hat{\beta}_i(\omega) = \beta_i\}.$$

As  $\hat{\beta}_i^l(\omega)(s_{-i}) = p^{m_l}([s_{-i}] \mid I_i(\omega))$ , it is the case that

$$\text{marg}_{S_j} \hat{\beta}_i^l(\omega)(s_j) = p^{m_l}([s_j] \mid I_i(\omega))$$

for all  $\omega \in \Omega$ , for all  $l = 1, \dots, L$ , for all  $s_j \in S_j$ , and for all  $j \in I \setminus \{i\}$ .

Epistemic hypotheses can be formalized by means of events. Some assumptions that will be used for the purpose of describing the interactive thinking underlying perfect equilibrium are now spelled out. The set

$$T_i := \{\omega \in \Omega : \hat{\beta}_i(\omega) \text{ is cautious}\}$$

denotes the event that *player  $i$  is cautious* and the event that *all players are cautious* is given by

$$T := \bigcap_{i \in I} T_i.$$

The set

$$R_i := \{\omega \in \Omega : \hat{s}_i(\omega) \text{ is lex-optimal given } \hat{\beta}_i(\omega)\}$$

constitutes the event that *player  $i$  is rational* and the event that *all players are rational* is denoted by

$$R := \bigcap_{i \in I} R_i.$$

Given some event  $E \subseteq \Omega$ , the set

$$PB_i(E) := \{\omega \in \Omega : p^{m_1}(E \mid I_i(\omega)) = 1\}$$

represents the event that *player  $i$  primarily believes in  $E$*  and the event that all players *primarily believe in  $E$*  is given by

$$PB := \bigcap_{i \in I} PB_i.$$

Note that primary belief concerns the first lexicographic *posterior* level  $l = m_1$  which may differ from the first lexicographic *prior* level  $m = 1$ .

As a preliminary observation we provide an epistemic foundation for perfect equilibrium in the special case of two player games.

**Proposition 2.** Let  $\Gamma$  be a finite game with two players  $i$  and  $j$ ,  $\mathcal{A}_{LCP}^{\Gamma}$  be some lexicographic Aumann model of  $\Gamma$ , and  $\omega^* \in \Omega$  be some world. If  $\omega^* \in PB(T) \cap PB(R) \cap K([\hat{\beta}_i(\omega^*)] \cap [\hat{\beta}_j(\omega^*)])$ , then there exists a pair of mixed strategies  $(\sigma_i, \sigma_j) \in \Delta(S_i) \times \Delta(S_j)$  such that

- $\sigma_i = \hat{b}_i^1(\omega^*)$  and  $\sigma_j = \hat{b}_j^1(\omega^*)$ ;
- the pair of mixed strategies  $(\sigma_i, \sigma_j)$  constitutes a perfect equilibrium of  $\Gamma$ .

**Proof.** (i) Define  $\beta_i := \hat{\beta}_i(\omega^*)$  and  $\beta_j := \hat{\beta}_j(\omega^*)$  as well as  $\sigma_i := b_i^1$  and  $\sigma_j := b_j^1$ . Then,  $\sigma_i = \hat{b}_i^1(\omega^*)$  and  $\sigma_j = \hat{b}_j^1(\omega^*)$  directly obtains.

(ii) Let  $k \in \{i, j\}$  be one of the two players and  $-k$  be his opponent. As  $\omega^* \in K([\hat{\beta}_i(\omega^*)] \cap [\hat{\beta}_j(\omega^*)]) \subseteq K_{-k}([\hat{\beta}_k(\omega^*)])$ , it follows that  $I_{-k}(\omega^*) \subseteq [\hat{\beta}_k(\omega^*)]$  and consequently  $\hat{\beta}_k(\omega) = \hat{\beta}_k(\omega^*)$  for all  $\omega \in I_{-k}(\omega^*)$ . As  $\omega^* \in PB(T) \subseteq PB_{-k}(T_k)$ , it is the case that

$$p^{m_1}(T_k \mid I_{-k}(\omega^*)) = \frac{p^{m_1}(T_k \cap I_{-k}(\omega^*))}{p^{m_1}(I_{-k}(\omega^*))} = 1$$

and thus there exists  $\omega' \in T_k \cap I_{-k}(\omega^*)$ . Then,  $\hat{\beta}_k(\omega')$  is cautious and  $\hat{\beta}_k(\omega') = \hat{\beta}_k(\omega^*) = \beta_k$ . It follows that  $\beta_k$  is cautious too. Since  $k$  has been chosen arbitrarily, property (a) of Definition 7 obtains. In addition,  $\sigma_k = b_{-k}^1 = \text{marg}_{S_k} b_{-k}^1$  ensures that property (b) of Definition 7 is satisfied. Next consider some strategy  $s_k \in \text{supp}(\sigma_k) = \text{supp}(\hat{b}_{-k}^1(\omega^*))$ . Then,  $\hat{b}_{-k}^1(\omega^*)(s_k) = p^{m_1}([s_k] \mid I_{-k}(\omega^*)) > 0$ , and thus there exists  $\omega' \in [s_k] \cap \text{supp}(p^{m_1}(\cdot \mid I_{-k}(\omega^*))) \subseteq [s_k] \cap I_{-k}(\omega^*)$ . Consequently,  $\hat{s}_k(\omega') = s_k$  and  $\hat{\beta}_k(\omega') = \hat{\beta}_k(\omega^*)$ . Also, as  $\omega^* \in PB(R) \subseteq PB_{-k}(R_k)$ , it is the case that  $p^{m_1}(R_k \mid I_{-k}(\omega^*)) = 1$  and thus  $\text{supp}(p^{m_1}(\cdot \mid I_{-k}(\omega^*))) \subseteq R_k$ . Hence,  $\omega' \in R_k$ , i.e.  $\hat{s}_k(\omega') = s_k$  is lex-optimal given  $\hat{\beta}_k(\omega') = \hat{\beta}_k(\omega^*) = \beta_k$ . This establishes property (c) of Definition 7. Besides, note that  $\beta_k = \text{marg}_{S_{-k}} \beta_k = \bigotimes_{j \in I \setminus \{k\}} \text{marg}_{S_j} \beta_k$  holds trivially as there is only one opponent for each player, which establishes property (d) of Definition 7. Finally, define  $\pi := \beta_i \otimes \beta_j$ . Then,  $\text{marg}_{S_{-k}} \pi = \beta_k$  directly follows, and property (e) of Definition 7 is satisfied. ■

The reasoning assumptions underlying perfect equilibrium, if attention is restricted to two players thus consist of mutual primary belief in caution, mutual primary belief in rationality, and mutual knowledge of conjectures.

In order to tame the complications arising once more than two players are admitted, the epistemic conditions need to be tightened. The problem of projection can be tackled by strengthening mutual knowledge of conjectures to common knowledge. By the aid of WAT, an epistemic foundation then ensues for the notion of lexicographic semi-perfect equilibrium.

**Lemma 2.** Let  $\Gamma$  be a finite game,  $\mathcal{A}_{LCP}^\Gamma$  be some lexicographic Aumann model of  $\Gamma$ , and  $\omega^* \in \Omega$  be some world. If  $\omega^* \in PB(T) \cap PB(R) \cap CK(\bigcap_{i \in I} [\hat{\beta}_i(\omega^*)])$ , then there exists a tuple of mixed strategies  $(\sigma_i^*)_{i \in I} \in \times_{i \in I} (\Delta(S_i))$  such that

- (i)  $\sigma_i^* = \arg_{S_i} \hat{b}_j^1(\omega^*)$  for all  $i \in I$  and for all  $j \in I \setminus \{i\}$ ;
- (ii) the tuple of mixed strategies  $(\sigma_i^*)_{i \in I}$  constitutes a lexicographic semi-perfect equilibrium of  $\Gamma$ .

**Proof.** (i) Consider the tuple of lexicographic conjectures  $(\hat{\beta}_i(\omega^*))_{i \in I}$  at world  $\omega^*$ . Let  $i \in I$  be some player. Observe that  $[\hat{\beta}_j(\omega^*)] \subseteq [\hat{b}_j^1(\omega^*)] \subseteq [\arg_{S_i}(\hat{b}_j^1(\omega^*))]$  for all  $j = 1, \dots, L$  and for all  $j \in I \setminus \{i\}$ . Then, by monotonicity of common knowledge,

$$CK\left(\bigcap_{j \in I \setminus \{i\}} [\hat{\beta}_j(\omega^*)]\right) \subseteq CK\left(\bigcap_{j \in I \setminus \{i\}} \bigcap_{l \in \{1, \dots, L\}} [\arg_{S_i} \hat{b}_j^l(\omega^*)]\right) \neq \emptyset.$$

As  $\arg_{S_i} \hat{b}_j^l(\omega^*)(s_i) = p^{m_l}([s_i] \mid I_j(\omega^*))$  for all  $\omega \in \Omega$ , for all  $j \in I \setminus \{i\}$ , for all  $l \in \{1, \dots, L\}$ , and for all  $s_i \in S_i$ ,

$$CK\left(\bigcap_{j \in I \setminus \{i\}} \bigcap_{l \in \{1, \dots, L\}} \{\omega \in \Omega : p^{m_l}([s_i] \mid I_j(\omega)) = p^{m_l}([s_i] \mid I_j(\omega^*))\}\right) \neq \emptyset$$

holds for all  $s_i \in S_i$ . By Theorem 1, it follows that

$$p^{m_l}([s_i] \mid I_j(\omega^*)) = p^{m_l}([s_i] \mid I_k(\omega^*))$$

for all  $s_i \in S_i$  as well as for all  $j, k \in I \setminus \{i\}$ , and thus

$$\arg_{S_i} \hat{b}_j^1(\omega^*) = \arg_{S_i} \hat{b}_k^1(\omega^*)$$

for all  $j, k \in I \setminus \{i\}$ . For every player  $i \in I$ , define  $\sigma_i^* := \arg_{S_i} \hat{b}_{i'}^1(\omega^*)$  for some  $i' \in I \setminus \{i\}$ . Then,  $\sigma_i^* = \arg_{S_i} \hat{b}_j^1(\omega^*)$  holds for all  $i \in I$  and for all  $j \in I \setminus \{i\}$ .

(ii) Consider the tuple of lexicographic conjectures  $(\hat{\beta}_i(\omega^*))_{i \in I}$ , where  $\hat{\beta}_i(\omega^*) = (\hat{b}_i^1(\omega^*), \dots, \hat{b}_i^L(\omega^*))$  for all  $i \in I$ . Let  $i, j \in I$  be two players such that  $i \neq j$ . Since  $\omega^* \in CK(\bigcap_{i \in I} [\hat{\beta}_i(\omega^*)]) \subseteq K_j([\hat{\beta}_i(\omega^*)])$ , it follows that  $I_j(\omega^*) \subseteq [\hat{\beta}_i(\omega^*)]$ , and thus  $\hat{\beta}_i(\omega) = \hat{\beta}_i(\omega^*)$  for all  $\omega \in I_j(\omega^*)$ . Note that  $\text{supp}(p^{m_1}(\cdot \mid I_j(\omega^*))) \subseteq I_j(\omega^*)$ . Moreover, as  $\omega^* \in PB_j(T_i)$ , the equation  $p^{m_1}(T_i \mid I_j(\omega^*)) = 1$  holds, thus  $\text{supp}(p^{m_1}(\cdot \mid I_j(\omega^*))) \subseteq T_i$ . Now, consider  $\omega' \in \text{supp}(p^{m_1}(\cdot \mid I_j(\omega^*)))$ . Then,  $\omega' \in I_j(\omega^*) \cap T_i$ . Consequently, on the one hand  $\hat{\beta}_i(\omega') = \hat{\beta}_i(\omega^*)$  and on the other hand  $\hat{\beta}_i(\omega')$  is cautious. Therefore,  $\hat{\beta}_i(\omega^*)$  is cautious, which establishes property (a) of Definition 8.

By part (i), the property  $\sigma_i^* = \arg_{S_i} \hat{b}_j^1(\omega^*)$  holds for all  $i \in I$  and for all  $j \in I \setminus \{i\}$ . Thus property (b) of Definition 8 obtains.

Let  $i, j \in I$  such that  $i \neq j$  and consider some  $s_i \in \text{supp}(\sigma_i^*) = \text{supp}(\arg_{S_i} \hat{b}_j^1(\omega^*))$ . Thus,  $\arg_{S_i} \hat{b}_j^1(\omega^*)(s_i) = p^{m_1}([s_i] \mid I_j(\omega^*)) > 0$ . Hence, there exists  $\omega^\circ \in \text{supp}(p^{m_1}(\cdot \mid I_j(\omega^*))) \subseteq I_j(\omega^*)$  such that  $\hat{s}_i(\omega^\circ) = s_i$ . As shown above, it is also the case that  $\hat{\beta}_i(\omega) = \hat{\beta}_i(\omega^*)$  for all  $\omega \in I_j(\omega^*)$ . Consequently,  $\hat{\beta}_i(\omega^\circ) = \hat{\beta}_i(\omega^*)$ . Since  $\omega^* \in PB(R) \subseteq PB_j(R_i)$ , it holds that  $p^{m_1}(R_i \mid I_j(\omega^*)) = 1$ , i.e.  $\omega' \in R_i$  for all  $\omega' \in \text{supp}(p^{m_1}(\cdot \mid I_j(\omega^*)))$ . Thus,  $\omega^\circ \in R_i$ , i.e.  $\hat{s}_i(\omega^\circ)$  is lex-optimal given  $\hat{\beta}_i(\omega^\circ)$ . As  $\hat{s}_i(\omega^\circ) = s_i$  and  $\hat{\beta}_i(\omega^\circ) = \hat{\beta}_i(\omega^*)$ , it follows that  $s_i$  is lex-optimal given  $\hat{\beta}_i(\omega^*)$ , which establishes property (c) of Definition 8.

Therefore,  $(\sigma_i^*)_{i \in I}$  constitutes a lexicographic semi-perfect equilibrium of  $\Gamma$ . ■

The weak agreement theorem (WAT) plays a major role in the preceding result, as it ensures that players always agree on their marginal conjectures about any common opponent they face in the game. The

possibility that any two players entertain distinct beliefs about a third player's choice is thereby blocked and the problem of projection solved. Formally, condition (i) of Lemma 2 and property (b) of Definition 8 are driven by WAT.

Yet additional armoury has to be invoked to establish an epistemic foundation for perfect equilibrium in the general set-up of many player games. Requiring the common prior to be mutually absolutely continuous enables the application of SAT, which can be used in turn to resolve the problem of independence.

**Theorem 3.** Let  $\Gamma$  be a finite game,  $\mathcal{A}_{LCP}^\Gamma$  be some lexicographic Aumann model of  $\Gamma$  such that the common prior  $\rho$  is mutually absolutely continuous, and  $\omega^* \in \Omega$  be some world. If  $\omega^* \in PB(T) \cap PB(R) \cap CK(\bigcap_{i \in I} [\hat{\beta}_i(\omega^*)])$ , then there exists a tuple of mixed strategies  $(\sigma_i^*)_{i \in I} \in \times_{i \in I} (\Delta(S_i))$  such that

- (i)  $\sigma_i^* = \arg_{S_i} \hat{b}_j^1(\omega^*)$  for all  $i \in I$  and for all  $j \in I \setminus \{i\}$ ;
- (ii) the tuple of mixed strategies  $(\sigma_i^*)_{i \in I}$  constitutes a perfect equilibrium of  $\Gamma$ .

**Proof.** (i) Consider the tuple of lexicographic conjectures  $(\hat{\beta}_i(\omega^*))_{i \in I}$  at world  $\omega^*$ . For every player  $i \in I$ , define  $\sigma_i^* := \arg_{S_i} \hat{b}_{i'}^1(\omega^*)$  for some  $i' \in I \setminus \{i\}$ . Part (i) of Lemma 2 ensures that  $\sigma_i^* = \arg_{S_i} \hat{b}_j^1(\omega^*)$  for all  $i \in I$  and for all  $j \in I \setminus \{i\}$ .

(ii) By Lemma 2, properties (a), (b), and (c) of Definition 7 hold. Let  $i \in I$  be some player and  $l \in \{1, \dots, L\}$  be some lexicographic level. Since  $CK(\bigcap_{j \in I} [\hat{\beta}_j(\omega^*)]) \neq \emptyset$ , it is the case that

$$CK([\arg_{S_i} \hat{b}_i^l(\omega^*)]) \neq \emptyset$$

$$CK([\arg_{S_{i+1}} \hat{b}_i^l(\omega^*)]) \neq \emptyset$$

$$CK\left(\bigcap_{j \in \{i, i+1\}} [\arg_{S_{-(i,j+1)}} \hat{b}_j^l(\omega^*)]\right) \neq \emptyset.$$

Consider some opponents' strategy combination  $s_{-i} \in S_{-i}$ . As  $\hat{b}_i^l(\omega^*)(\cdot) = p^{m_l}(\cdot \mid I_i(\omega^*))$ , it follows that

$$CK\left(\{\omega \in \Omega : p^{m_l}([s_{-i}] \mid I_i(\omega)) = p^{m_l}([s_{-i}] \mid I_i(\omega^*))\}\right) \neq \emptyset$$

$$CK\left(\{\omega \in \Omega : p^{m_l}([s_{i+1}] \mid I_i(\omega)) = p^{m_l}([s_{i+1}] \mid I_i(\omega^*))\}\right) \neq \emptyset$$

$$CK\left(\bigcap_{j \in \{i, i+1\}} \{\omega \in \Omega : p^{m_l}([s_{-(i,j+1)}] \mid I_j(\omega)) = p^{m_l}([s_{-(i,j+1)}] \mid I_j(\omega^*))\}\right) \neq \emptyset$$

By the proof of Theorem 2, there exist some indices  $\alpha_i$ ,  $\beta_i$  and  $\gamma_i$  independent from  $i$ ,  $i+1$  and  $\omega$  such that

$$p^{m_l}([s_{-i}] \mid I_i(\omega)) = p^{\alpha_i}([s_{-i}] \mid (\bigwedge_{i' \in I} I_{i'})(\omega^*))$$

$$p^{m_l}([s_{i+1}] \mid I_i(\omega)) = p^{\beta_i}([s_{i+1}] \mid (\bigwedge_{i' \in I} I_{i'})(\omega^*))$$

$$p^{m_l}([s_{-(i,j+1)}] \mid I_j(\omega)) = p^{m_l}([s_{-(i,j+1)}] \mid I_{i+1}(\omega)) = p^{\gamma_i}([s_{-(i,j+1)}] \mid (\bigwedge_{i' \in I} I_{i'})(\omega^*))$$

for all  $\omega \in (\bigwedge_{i \in I} I_i)(\omega^*)$ . Since  $\rho$  is mutually absolutely continuous, the first part of the proof of Theorem 2 ensures that the lexicographic levels of  $\rho(\cdot \mid I_i(\omega))$  are the same for all  $\omega \in (\bigwedge_{i \in I} I_i)(\omega^*)$ , and thus  $\alpha_i = \beta_i = \gamma_i := \bar{m}_l$ . Moreover, since either  $I_i(\omega) \subseteq [s_i]$  or  $I_i(\omega) \subseteq [s_i]^c$ , the following property holds

$$p(E \cap [s_i] \mid I_i(\omega)) = p(E \mid I_i(\omega)) \cdot p([s_i] \mid I_i(\omega))$$

for all probability measures  $p \in \Delta(\Omega)$ , for all  $E \subseteq \Omega$  and for all  $i \in I$ . Let  $\mathcal{P} := \{P_{i+1} \in I_{i+1} : P_{i+1} \subseteq (\bigwedge_{i \in I} I_i)(\omega^*)\}$  be the possibility cells of player  $i+1$  included in the meet cell of  $\omega^*$ . By using the above properties together with the law of total probability, it follows that

$$\begin{aligned} & p^{m_l}([s_{-i}] \mid I_i(\omega^*)) \\ &= p^{\bar{m}_l}([s_{-i}] \mid (\bigwedge_{j \in I} I_j)(\omega^*)) \\ &= \sum_{P_{i+1} \in \mathcal{P}} p^{\bar{m}_l}([s_{-i}] \mid P_{i+1}) \cdot p^{\bar{m}_l}(P_{i+1} \mid (\bigwedge_{j \in I} I_j)(\omega^*)) \end{aligned}$$

$$\begin{aligned}
&= \sum_{P_{i+1} \in \mathcal{P}} p^{\bar{m}_i}([s_{-(i,i+1)}] | P_{i+1}) \cdot p^{\bar{m}_i}([s_{i+1}] | P_{i+1}) \cdot p^{\bar{m}_i}(P_{i+1} | (\bigwedge_{j \in I} I_j)(\omega^*)) \\
&= \sum_{P_{i+1} \in \mathcal{P}} p^{\bar{m}_i}([s_{-(i,i+1)}] | (\bigwedge_{j \in I} I_j)(\omega^*)) \cdot p^{\bar{m}_i}([s_{i+1}] | P_{i+1}) \\
&\quad \cdot p^{\bar{m}_i}(P_{i+1} | (\bigwedge_{j \in I} I_j)(\omega^*)) \\
&= p^{\bar{m}_i}([s_{-(i,i+1)}] | (\bigwedge_{j \in I} I_j)(\omega^*)) \cdot \sum_{P_{i+1} \in \mathcal{P}} p^{\bar{m}_i}([s_{i+1}] | P_{i+1}) \\
&\quad \cdot p^{\bar{m}_i}(P_{i+1} | (\bigwedge_{j \in I} I_j)(\omega^*)) \\
&= p^{\bar{m}_i}([s_{-(i,i+1)}] | (\bigwedge_{j \in I} I_j)(\omega^*)) \cdot p^{\bar{m}_i}([s_{i+1}] | (\bigwedge_{j \in I} I_j)(\omega^*)) \\
&= p^{\bar{m}_i}([s_{-(i,i+1)}] | I_{i+1}(\omega^*)) \cdot p^{\bar{m}_i}([s_{i+1}] | I_i(\omega^*)).
\end{aligned}$$

Analogously,

$$p^{\bar{m}_i}([s_{-(i,i+1)}] | I_{i+1}(\omega^*)) = p^{\bar{m}_i}([s_{-(i,i+1,i+2)}] | I_{i+2}(\omega^*)) \cdot p^{\bar{m}_i}([s_{i+2}] | I_{i+1}(\omega^*))$$

ensues, and thus

$$p^{\bar{m}_i}([s_{-i}] | I_i(\omega^*)) = p^{\bar{m}_i}([s_{-(i,i+1,i+2)}] | I_{i+2}(\omega^*)) \cdot p^{\bar{m}_i}([s_{i+2}] | I_{i+1}(\omega^*)) \cdot p^{\bar{m}_i}([s_{i+1}] | I_i(\omega^*)).$$

By induction, it follows that

$$p^{\bar{m}_i}([s_{-i}] | I_i(\omega^*)) = \prod_{j \in I \setminus \{i-1\}} p^{\bar{m}_i}([s_{j+1}] | I_j(\omega^*)).$$

Consequently,

$$\begin{aligned}
&\hat{b}_i^l(\omega^*)(s_{-i}) \\
&= p^{\bar{m}_i}([s_{-i}] | I_i(\omega^*)) \\
&= \prod_{j \in I \setminus \{i-1\}} p^{\bar{m}_i}([s_{j+1}] | I_j(\omega^*)) = \prod_{j \in I \setminus \{i-1\}} \text{marg}_{S_{j+1}} \hat{b}_j^l(\omega^*)(s_{j+1}) \\
&= \prod_{j \in I \setminus \{i-1\}} \text{marg}_{S_{j+1}} \hat{b}_i^l(\omega^*)(s_{j+1}) = \prod_{j \in I \setminus \{i\}} \text{marg}_{S_j} \hat{b}_i^l(\omega^*)(s_j).
\end{aligned}$$

Therefore,  $\hat{b}_i(\omega^*) = \bigotimes_{j \in I \setminus \{i\}} \text{marg}_{S_j} \hat{b}_i(\omega^*)$ , which establishes property (d) of Definition 7.

Furthermore, let  $i \in I$  and  $j, j' \in I \setminus \{i\}$  be some players,  $s_i \in S_i$  be some strategy for player  $i$ , and  $l \in \{1, \dots, L\}$  be some lexicographic level. Observe that

$$\begin{aligned}
\text{marg}_{S_i} \hat{b}_j^l(\omega^*)(s_i) &= p^{\bar{m}_i}([s_i] | I_j(\omega^*)) = p^{\bar{m}_i}([s_i] | (\bigwedge_{i' \in I} I_{i'})(\omega^*)) \\
&= p^{\bar{m}_i}([s_i] | I_{j'}(\omega^*)) = \text{marg}_{S_{i'}} \hat{b}_{j'}^l(\omega^*)(s_i)
\end{aligned}$$

and therefore,  $\text{marg}_{S_i} \hat{b}_j(\omega^*) = \text{marg}_{S_i} \hat{b}_{j'}(\omega^*)$  for all  $i \in I$  and for all  $j, j' \in I \setminus \{i\}$ . Now, take  $i, i' \in I$  such that  $i \neq i'$  and define the lexicographic product measure

$$\pi := \hat{b}_i(\omega^*) \otimes \text{marg}_{S_{i'}} \hat{b}_{i'}(\omega^*) = \left( \bigotimes_{j \in I \setminus \{i\}} \text{marg}_{S_j} \hat{b}_i(\omega^*) \right) \otimes \text{marg}_{S_{i'}} \hat{b}_{i'}(\omega^*).$$

We show that  $\text{marg}_{S_k} \pi = \hat{b}_k(\omega^*)$  for all  $k \in I$ . First, the definition of  $\pi$  combined with property (d) of Definition 7 ensures that

$$\text{marg}_{S_{-i}} \pi = \bigotimes_{j \in I \setminus \{i\}} \text{marg}_{S_j} \hat{b}_i(\omega^*) = \hat{b}_i(\omega^*).$$

If  $k \in I \setminus \{i\}$ , then the equality of the marginal conjectures established above together with property (d) of Definition 7 implies that

$$\begin{aligned}
\text{marg}_{S_k} \pi &= \left( \bigotimes_{j \in I \setminus \{i,k\}} \text{marg}_{S_j} \hat{b}_i(\omega^*) \right) \otimes \text{marg}_{S_{i'}} \hat{b}_{i'}(\omega^*) \\
&= \left( \bigotimes_{j \in I \setminus \{i,k\}} \text{marg}_{S_j} \hat{b}_k(\omega^*) \right) \otimes \text{marg}_{S_{i'}} \hat{b}_k(\omega^*) \\
&= \bigotimes_{j \in I \setminus \{k\}} \text{marg}_{S_j} \hat{b}_k(\omega^*) = \hat{b}_k(\omega^*).
\end{aligned}$$

Consequently,  $\pi$  and  $(\hat{b}_i(\omega^*))_{i \in I}$  satisfy property (e) of Definition 7.

Therefore,  $(\sigma_i^*)_{i \in I}$  forms a lexicographic perfect equilibrium of  $\Gamma$ , and thus, by Lemma 1, a perfect equilibrium of  $\Gamma$ . ■

The property that a player's belief about his opponents' strategies is independent poses a rather intricate matter in the proof of Theorem 3 and its accomplishment is assisted by our strong agreement theorem (SAT). The effective application of the two lexicographic agreement theorems (WAT and SAT) in establishing epistemic conditions for perfect equilibrium once again underlines the power that Aumann's seminal impossibility result on agreeing to disagree is capable of unfolding.

The following result addresses the converse direction by ensuring that the epistemic conditions of Theorem 3 always exist and can be aligned with any perfect equilibrium.

**Theorem 4.** Let  $\Gamma$  be a finite game and  $\sigma = (\sigma_i)_{i \in I} \in \times_{i \in I} (\Delta(S_i))$  be a tuple of mixed strategies that constitutes a perfect equilibrium of  $\Gamma$ . Then, there exists a lexicographic Aumann model  $\mathcal{A}_{LCP}^{\Gamma}$  of  $\Gamma$  with a world  $\omega^* \in \Omega$  such that the common prior  $\rho$  is mutually absolutely continuous,  $\omega^* \in PB(T) \cap PB(R) \cap CK(\bigcap_{i \in I} [\hat{b}_i(\omega^*)])$ , as well as  $\sigma_i = \text{marg}_{S_i} \hat{b}_i^1(\omega^*)$  for all  $i \in I$  and for all  $j \in I \setminus \{i\}$ .

**Proof.** By Lemma 1,  $\sigma$  forms a lexicographic perfect equilibrium and there exist a tuple  $\beta = (\beta_i)_{i \in I} \in \left( (\Delta(S_{-i}))^L \right)_{i \in I}$  of lexicographic conjectures and a lexicographic product measure  $\pi = (\pi^1, \dots, \pi^L) \in (\Delta(\times_{i \in I} S_i))^L$  in line with the properties (a) to (e) of Definition 7. Construct the lexicographic Aumann model  $\mathcal{A}_{LCP}^{\Gamma} = (\Omega, \rho, I, (I_i, \hat{s}_i)_{i \in I})$  of  $\Gamma$ , where

- $\Omega = \{\omega^s : s = (s_i)_{i \in I} \in \times_{i \in I} S_i\}$ ,
- $\rho^m \in \Delta(\Omega)$  is defined by  $\rho^m(\omega^s) = \pi^m(s)$  for all  $\omega^s \in \Omega$  and for all  $m \in \{1, \dots, M\}$ , with  $M = L$ ,
- $I_i(\omega^s) = \Omega$  for all  $\omega^s \in \Omega$  and for all  $i \in I$ ,
- $\hat{s}_i : \Omega \rightarrow \times_{i \in I} S_i$  is defined by  $\hat{s}_i(\omega^s) = s_i$ , for all  $\omega^s \in \Omega$  and for all  $i \in I$ .

As  $I_i(\omega^s) = \Omega$  for all  $\omega^s \in \Omega$  and for all  $i \in I$ , it directly follows that  $\rho^m(I_i(\omega^s)) = \rho^m(I_j(\omega^s)) = 1$ , and thus  $\rho^l(I_i(\omega^s)) = 0$  if and only if  $\rho^l(I_j(\omega^s)) = 0$ , for all  $\omega^s \in \Omega$ , for all  $m \in \{1, \dots, M\}$ , and for all  $i, j \in I$ . Therefore,  $\rho$  is mutually absolutely continuous.

Since  $\rho^m(I_i(\omega^s)) = 1$  for all  $\omega^s \in \Omega$ , for all  $m \in \{1, \dots, M\}$ , and for all  $i \in I$ , Definition 3 ensures that  $m_l = l$  for all  $l \in \{1, \dots, L\}$ . Consider some player  $i \in I$ , some world  $\omega^s \in \Omega$ , and some lexicographic level  $l \in \{1, \dots, L\}$ . It follows that

$$\begin{aligned}
\hat{b}_i^l(\omega^s)(s'_{-i}) &= p^{\bar{m}_i}([s'_{-i}] | I_i(\omega^s)) = \frac{p^l(I_i(\omega^s) \cap [s'_{-i}])}{p^l(I_i(\omega^s))} \\
&= \frac{p^l(\Omega \cap [s'_{-i}])}{p^l(\Omega)} = \frac{p^l([s'_{-i}])}{1} \\
&= \pi^l(\{s \in \times_{i \in I} S_i : s_{-i} = s'_{-i}\}) = \sum_{s_i \in S_i} \pi^l(s_i, s'_{-i}) \\
&= \text{marg}_{S_{-i}} \pi^l(s'_{-i}) = \hat{b}_i^l(s'_{-i})
\end{aligned}$$

for all  $s'_{-i} \in S_{-i}$ , where the last equality is due to property (e) of Definition 7. Consequently,  $\hat{b}_i(\omega^s) = \beta_i$  for all  $\omega^s \in \Omega$  and for all  $i \in I$ . Hence,  $[\hat{b}_i(\omega^s)] = \{\omega^s \in \Omega : \hat{b}_i(\omega^s) = \beta_i(\omega^s)\} = \{\omega^s \in \Omega : \hat{b}_i(\omega^s) = \beta_i\} = \Omega$  for all  $\omega^s \in \Omega$  as well as for all  $i \in I$ , and thus  $CK(\bigcap_{i \in I} [\hat{b}_i(\omega^s)]) = CK(\Omega) = \Omega$ .

Next consider some world  $\omega^s \in \Omega$  and some player  $i \in I$ . Since  $\hat{b}_i(\omega^s) = \beta_i$ , property (a) of lexicographic perfect equilibrium ensures that  $\hat{b}_i(\omega^s)$  is cautious, i.e.  $\omega^s \in T_i$ . It follows that  $T_i = \Omega$ , and thus  $T = \bigcap_{i \in I} T_j = \Omega$ . Consequently,  $\text{supp}(p^{\bar{m}_i}(\cdot | I_i(\omega^s))) \subseteq T$  and hence  $p^{\bar{m}_i}(T | I_i(\omega^s)) = 1$ , i.e.  $\omega^s \in PB_i(T)$ . Also, by properties (b) and (e) of Definition 7, it follows that

$$p^{\bar{m}_i}(\cdot | I_i(\omega^s)) = p^1(\cdot | \Omega) = p^1 = \pi^1 = \bigotimes_{j \in I} \text{marg}_{S_j} \pi^1$$

$$= \bigotimes_{j \in I} \text{marg}_{S_j} \text{marg}_{S_{-(j+1)}} \pi^1 = \bigotimes_{j \in I} \text{marg}_{S_j} b_{j+1}^1 = \bigotimes_{j \in I} \sigma_j$$

Let  $\omega^{s'} \in \text{supp}(p^{m_1}(\cdot \mid I_i(\omega^s)))$ . Then,  $s' \in \text{supp}(\bigotimes_{j \in I} \sigma_j)$ , i.e.  $s'_j \in \text{supp}(\sigma_j)$  for all  $j \in I$ . By property (c) of Definition 7,  $s'_j$  is lex-optimal given  $\beta_j$ , and hence  $\hat{s}_j(\omega^{s'})$  is lex-optimal given  $\hat{\beta}_j(\omega^{s'})$ , i.e.  $\omega^{s'} \in R_j$  for all  $j \in I$ . Thus  $\omega^{s'} \in \bigcap_{j \in I} R_j = R$ . Hence,  $\text{supp}(p^{m_1}(\cdot \mid I_i(\omega^s))) \subseteq R$ . Thus,  $p^{m_1}(R \mid I_i(\omega^s)) = 1$ , i.e.  $\omega^s \in PB_i(R)$ . Since  $i$  has been chosen arbitrarily,  $\omega^s \in \bigcap_{i \in I} PB_i(T) \cap \bigcap_{i \in I} PB_i(R) = PB(T) \cap PB(R)$ . As  $\omega^s$  has been picked arbitrarily too,  $PB(T) \cap PB(R) = \Omega$  obtains.

Finally, let  $\omega^* \in \Omega$  be some world and  $i \in I$  be some player. Then,  $\omega^* \in PB(T) \cap PB(R) \cap CK(\bigcap_{i \in I} [\hat{\beta}_i(\omega^*)])$ . Furthermore, property (b) of Definition 7 guarantees that  $\sigma_i = \text{marg}_{S_i} b_j^1 = \text{marg}_{S_i} \hat{b}_j^1(\omega^*)$  for all  $j \in I \setminus \{i\}$ . Since  $i$  has been chosen arbitrarily,  $\sigma_i = \text{marg}_{S_i} \hat{b}_j^1(\omega^*)$  for all  $i \in I$  and for all  $j \in I \setminus \{i\}$ . ■

Accordingly, the sufficient conditions for perfect equilibrium put forth by Theorem 3 are not too strong in the sense that every perfect equilibrium is attainable with them. The conjunction of Theorems 3 and 4 constitutes an epistemic characterization of perfect equilibrium in terms of mutual primary belief in caution, mutual primary belief in rationality, and common knowledge of conjectures.

The epistemic programme in game theory has shed light on the reasoning assumptions underlying Nash equilibrium.<sup>17</sup> The decisive – yet conceptually not unproblematic – implicit property of Nash equilibrium lies in some correct beliefs assumption. By requiring common knowledge of conjectures, Theorems 3 and 4 show that a significant dose of doxastic inerrancy also underlies the more general solution concept of perfect equilibrium. In contrast, common knowledge of rationality is not required in terms of reasoning: it is not even needed at the first lexicographic level. A central conceptual insight due to Aumann and Brandenburger (1995) for Nash equilibrium – interactive beliefs in rationality do not enter the picture but only an interactive correct beliefs condition does – is thus fortified by Theorems 3 and 4 in the more general context of perfect equilibrium.<sup>18</sup> Both Nash equilibrium and perfect equilibrium hence only require iterated – and thus truly interactive – beliefs about conjectures and not about rationality or anything else. Consequently, some correct beliefs property constitutes the essence of these solution concepts. Nonetheless, the reasoning foundation for perfect equilibrium stretches beyond the one for Nash equilibrium. Indeed, some notion of caution is needed in order to reflect the inherent trembles property of perfect equilibrium, which is absent from Nash equilibrium though.

## 9. Conclusion

When interactive epistemology is enriched by lexicographic probability systems, three results on agreeing to disagree obtain. If the agents' posteriors are common knowledge, the weak agreement theorem ensures the first lexicographic level posteriors to coincide. Somewhat unexpectedly, however, disagreement cannot be excluded without further assumptions on the deeper lexicographic levels. In line with our disagreement result, agreement can already fail on the second lexicographic level. Imposing mutual absolute continuity on top of common knowledge of posteriors, the strong agreement theorem rules out posterior disagreement at any lexicographic level.

<sup>17</sup> For instance, (Brandenburger, 1992a; Aumann and Brandenburger, 1995; Perea, 2007; Barelli, 2009; Bach and Tsakas, 2014; Bonanno, 2018; Bach and Perea, 2020).

<sup>18</sup> As Aumann and Brandenburger (1995) as well as Brandenburger (1992a) highlight, common knowledge enters the picture in an unexpected way for Nash equilibrium to ensue: what is needed is common knowledge of the players' conjectures but not of the players' rationality (cf. (Aumann and Brandenburger, 1995), p. 1163), and then only in games with more than two players (cf. (Brandenburger, 1992a), p. 96).

The impossibility of lexicographic agreeing to disagree becomes an essential tool to shed light on interactive reasoning in games. Epistemic conditions are provided for the classical solution concept of perfect equilibrium. In particular, the weak agreement theorem and the strong agreement theorem fundamentally assist in overcoming the challenges that arise with more than two players. The reasoning assumptions underlying perfect equilibrium are identified in our lexicographic framework by mutual primary belief in caution, mutual primary belief in rationality, and common knowledge of conjectures. The solution concept's key epistemic ingredient thus lies in an interactive correct beliefs assumption, while caution as well as rationality only appear in a non-iterated doxastic way on the first lexicographic level.

From a conceptual perspective, our results on the (im)possibility of lexicographic agreeing to disagree are relevant for situations when reasoning about ordered layers of contingencies is considered. Notably the original conclusion of Aumann's agreement theorem breaks down. Agreeing to disagree becomes conceivable once hypothetical contingencies enter the picture. This could have intriguing consequences for economic applications such as the possibility of trade. We leave such considerations for future research.

## Declaration of competing interest

We, the authors, declare that we have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper entitled "Lexicographic Agreeing to Disagree and Perfect Equilibrium".

## Data availability

No data was used for the research described in the article.

## Appendix

The proof of Lemma 1 requires some additional results that are laid out first.

Given a game  $\Gamma$ , some player  $i \in I$ , some strategy  $s_i \in S_i$  of player  $i$ , and some general – not necessarily product – probability measure  $q \in \Delta(S_{-i})$ , player  $i$ 's expected utility of strategy  $s_i$  is defined as

$$u_i(s_i, q) := \sum_{s_{-i} \in S_{-i}} q(s_{-i}) \cdot U_i(s_i, s_{-i}).$$

Let  $X$  be a finite space, let  $L > 0$  be an integer, let  $\alpha = (a^1, \dots, a^L) \in (\Delta(X))^L$  be a tuple of probability measures and let  $r = (r^1, \dots, r^{L-1}) \in (0, 1)^{L-1}$  be a tuple of real numbers. Let  $r \sqcap \alpha$  be defined by

$$r \sqcap \alpha := \begin{cases} a^1 & \text{if } L = 1 \\ (1 - r^1) \cdot a^1 + r^1 \cdot (1 - r^2) \cdot a^2 + r^1 \cdot r^2 \cdot (1 - r^3) \cdot a^3 + \dots + r^1 \cdot r^2 \cdot \dots \cdot r^{L-2} \cdot (1 - r^{L-1}) \cdot a^{L-1} + r^1 \cdot r^2 \cdot \dots \cdot r^{L-1} \cdot a^L & \text{if } L > 1 \end{cases}$$

Observe that  $r \sqcap \alpha \in \Delta(X)$ , since

$$\sum_{x \in X} (r \sqcap \alpha)(x) = (1 - r^1) + r^1 \cdot (1 - r^2) + r^1 \cdot r^2 \cdot (1 - r^3) + \dots + r^1 \cdot r^2 \cdot \dots \cdot r^{L-2} \cdot (1 - r^{L-1}) + r^1 \cdot r^2 \cdot \dots \cdot r^{L-1} = 1.$$

**Lemma A.1.** Let  $\Gamma$  be a game,  $i \in I$  be a player,  $s'_i, s''_i \in S_i$  be two strategies of player  $i$ ,  $\beta_i = (b^1, \dots, b^L) \in (\Delta(S_{-i}))^L$  be a lexicographic conjecture of player  $i$ , and  $(r_n)_{n \in \mathbb{N}} = ((r_n^1, \dots, r_n^{L-1}))_{n \in \mathbb{N}} \in [(0, 1)^{L-1}]^{\mathbb{N}}$  be a sequence such that  $\lim_{n \rightarrow \infty} r_n = \vec{0} \in \mathbb{R}^{L-1}$ . Then, the following properties hold:

- (i) If  $u_i(s'_i, r_n \sqcap \beta_i) > u_i(s''_i, r_n \sqcap \beta_i)$  for all  $n \in \mathbb{N}$ , then  $i$  strictly lex-prefers  $s'_i$  to  $s''_i$ .

- (ii) If  $u_i(s'_i, r_n \square \beta_i) \geq u_i(s_i, r_n \square \beta_i)$  for all  $n \in \mathbb{N}$  and for all  $s_i \in S_i$ , then  $s'_i$  is lex-optimal given  $\beta_i$ .
- (iii) If  $s'_i$  is lex-optimal given  $\beta_i$ , then there exist a subsequence  $(r_{n_k})_{k \in \mathbb{N}}$  of  $(r_n)_{n \in \mathbb{N}}$  and an index  $K \in \mathbb{N}$  such that  $u_i(s'_i, r_{n_k} \square \beta_i) \geq u_i(s_i, r_{n_k} \square \beta_i)$  for all  $k \geq K$  and for all  $s_i \in S_i$ .

**Proof.** (i) Observe that  $\lim_{n \rightarrow \infty} r_n = \vec{0}$  implies

$$\lim_{n \rightarrow \infty} r_n \square \beta_i = \lim_{n \rightarrow \infty} [(1-r_n^1) \cdot b_i^1 + r_n^1 \cdot (1-r_n^2) \cdot b_i^2 + \dots + r_n^1 \cdot r_n^2 \cdot \dots \cdot r_n^{L-1} \cdot b_i^L] = b_i^1.$$

In addition, for each  $l \in \{1, \dots, L\}$ , define

$$\Delta^l := u_i'(s'_i, \beta_i) - u_i'(s_i, \beta_i) = \sum_{s_{-i} \in S_{-i}} b_i^l(s_{-i}) \cdot [U_i(s'_i, s_{-i}) - U_i(s_i, s_{-i})].$$

Suppose that  $u_i(s'_i, r_n \square \beta_i) > u_i(s'_i, r_n \square \beta_i)$  for all  $n \in \mathbb{N}$ . It follows that

$$\begin{aligned} & \sum_{s_{-i} \in S_{-i}} (r_n \square \beta_i)(s_{-i}) \cdot [U_i(s'_i, s_{-i}) - U_i(s_i, s_{-i})] \\ &= (1-r_n^1) \cdot \Delta^1 + r_n^1 \cdot (1-r_n^2) \cdot \Delta^2 + \dots + r_n^1 \cdot r_n^2 \cdot \dots \cdot r_n^{L-1} \cdot \Delta^L > 0 \end{aligned} \quad (2)$$

for all  $n \in \mathbb{N}$ . Consequently,

$$\begin{aligned} 0 &\leq \lim_{n \rightarrow \infty} \sum_{s_{-i} \in S_{-i}} (r_n \square \beta_i)(s_{-i}) \cdot [U_i(s'_i, s_{-i}) - U_i(s_i, s_{-i})] \\ &= \sum_{s_{-i} \in S_{-i}} \lim_{n \rightarrow \infty} (r_n \square \beta_i)(s_{-i}) \cdot [U_i(s'_i, s_{-i}) - U_i(s_i, s_{-i})] \\ &= \sum_{s_{-i} \in S_{-i}} b_i^1(s_{-i}) \cdot [U_i(s'_i, s_{-i}) - U_i(s_i, s_{-i})] = \Delta^1. \end{aligned}$$

If  $\Delta^1 > 0$ , then  $u_i^1(s'_i, \beta_i) > u_i^1(s_i, \beta_i)$  and thus  $i$  strictly lex-prefers  $s'_i$  to  $s_i$ . If  $\Delta^1 = 0$ , then define the truncated tuples  $\beta_i^{(2)} := (b_i^2, \dots, b_i^L)$  and  $(r_n^{(2)})_{n \in \mathbb{N}} := ((r_n^2, \dots, r_n^{L-1}))_{n \in \mathbb{N}}$ . Property (2) together with the fact that  $\Delta^1 = 0$  ensures that

$$\begin{aligned} 0 &< \sum_{s_{-i} \in S_{-i}} (r_n \square \beta_i)(s_{-i}) \cdot [U_i(s'_i, s_{-i}) - U_i(s_i, s_{-i})] \\ &= (1-r_n^1) \cdot \Delta^1 + r_n^1 \cdot \sum_{s_{-i} \in S_{-i}} (r_n^{(2)} \square \beta_i^{(2)})(s_{-i}) \cdot [U_i(s'_i, s_{-i}) - U_i(s_i, s_{-i})] \\ &= r_n^1 \cdot \sum_{s_{-i} \in S_{-i}} (r_n^{(2)} \square \beta_i^{(2)})(s_{-i}) \cdot [U_i(s'_i, s_{-i}) - U_i(s_i, s_{-i})] \end{aligned}$$

for all  $n \in \mathbb{N}$ , and thus

$$\sum_{s_{-i} \in S_{-i}} (r_n^{(2)} \square \beta_i^{(2)})(s_{-i}) \cdot [U_i(s'_i, s_{-i}) - U_i(s_i, s_{-i})] > 0$$

for all  $n \in \mathbb{N}$ . Consequently,

$$\begin{aligned} 0 &\leq \lim_{n \rightarrow \infty} \sum_{s_{-i} \in S_{-i}} (r_n^{(2)} \square \beta_i^{(2)})(s_{-i}) \cdot [U_i(s'_i, s_{-i}) - U_i(s_i, s_{-i})] \\ &= \sum_{s_{-i} \in S_{-i}} \lim_{n \rightarrow \infty} (r_n^{(2)} \square \beta_i^{(2)})(s_{-i}) \cdot [U_i(s'_i, s_{-i}) - U_i(s_i, s_{-i})] \\ &= \sum_{s_{-i} \in S_{-i}} b_i^2(s_{-i}) \cdot [U_i(s'_i, s_{-i}) - U_i(s_i, s_{-i})] = \Delta^2. \end{aligned}$$

If  $\Delta^2 > 0$ , then  $u_i^2(s'_i, \beta_i) > u_i^2(s_i, \beta_i)$  and  $u_i^1(s'_i, \beta_i) = u_i^1(s_i, \beta_i)$ , and thus  $i$  strictly lex-prefers  $s'_i$  to  $s_i$ . If  $\Delta^2 = 0$ , then by continuing in this fashion for  $l \geq 3$ , property (2) ensures that eventually there exists  $l^* \in \{1, \dots, L\}$  such that  $\Delta^{l^*} > 0$  and  $\Delta^l = 0$  for all  $0 < l < l^*$ . Equivalently,  $u_i^{l^*}(s'_i, \beta_i) > u_i^{l^*}(s_i, \beta_i)$  and  $u_i^l(s'_i, \beta_i) = u_i^l(s_i, \beta_i)$  for all  $0 < l < l^*$ . Therefore,  $i$  strictly lex-prefers  $s'_i$  to  $s_i$ .

(ii) Let  $s_i \in S_i$ . Suppose that  $u_i(s'_i, r_n \square \beta_i) \geq u_i(s_i, r_n \square \beta_i)$  for all  $n \in \mathbb{N}$ . If  $u_i(s'_i, r_n \square \beta_i) = u_i(s_i, r_n \square \beta_i)$  for all  $n \in \mathbb{N}$ , then by similar arguments as in the proof of Lemma A.1 (i), it follows that  $\Delta^l = 0$  for all  $l \in \{1, \dots, L\}$ . Consequently,  $i$  weakly lex-prefers  $s'_i$  to  $s_i$ . If  $u_i(s'_i, r_n \square \beta_i) > u_i(s_i, r_n \square \beta_i)$  for some  $n^* \in \mathbb{N}$ , then again by similar arguments as in the proof of Lemma A.1 (i), there exists  $l^* \in \{1, \dots, L\}$  such that  $\Delta^{l^*} > 0$  and  $\Delta^l = 0$  for all  $0 < l < l^*$ . Hence,  $i$  weakly lex-prefers  $s'_i$  to  $s_i$  and, as  $s_i$  has been chosen arbitrarily,  $s'_i$  is thus lex-optimal given  $\beta_i$ .

(iii) Consider a subsequence  $(r_{n_k})_{k \in \mathbb{N}}$  of  $(r_n)_{n \in \mathbb{N}}$  that satisfies the following property: for every  $s_i \in S_i$ , if  $u_i(s_i, r_{n_k} \square \beta_i) > u_i(s'_i, r_{n_k} \square \beta_i)$

for infinitely many indices  $k \in \mathbb{N}$ , then  $u_i(s_i, r_{n_k} \square \beta_i) > u_i(s'_i, r_{n_k} \square \beta_i)$  all  $k \in \mathbb{N}$ . Since  $\lim_{n \rightarrow \infty} r_n = \vec{0}$ , it is the case that  $\lim_{k \rightarrow \infty} r_{n_k} = \vec{0}$ . Suppose that  $s'_i$  is lex-optimal given  $\beta_i$ . By the contraposition of Lemma A.1 (i), for all  $s_i \in S_i$ , it is not the case that  $u_i(s_i, r_{n_k} \square \beta_i) > u_i(s'_i, r_{n_k} \square \beta_i)$  for all  $k \in \mathbb{N}$ . The contraposition of the property of the sequence  $(r_{n_k})_{k \in \mathbb{N}}$  then ensures that, for all  $s_i \in S_i$ , it is not the case that  $u_i(s_i, r_{n_k} \square \beta_i) > u_i(s'_i, r_{n_k} \square \beta_i)$  for infinitely many indices  $k \in \mathbb{N}$ . Equivalently, for all  $s_i \in S_i$ , there exists  $K(s_i) \in \mathbb{N}$  such that  $u_i(s'_i, r_{n_k} \square \beta_i) \geq u_i(s_i, r_{n_k} \square \beta_i)$  for all  $k \geq K(s_i)$ . Consequently,  $u_i(s'_i, r_{n_k} \square \beta_i) \geq u_i(s_i, r_{n_k} \square \beta_i)$  for all  $k \geq \max\{K(s_i) : s_i \in S_i\}$  and for all  $s_i \in S_i$ . ■

**Lemma A.2.** Let  $\Gamma$  be a game and  $\psi : \Delta(\times_{i \in I} S_i) \times \Delta(\times_{i \in I} S_i) \rightarrow \mathbb{R}$  be the function defined by

$$\psi(\sigma, \tilde{\sigma}) := \sup\{r \in \mathbb{R} : \sigma(s) - r \cdot \tilde{\sigma}(s) \geq 0, \text{ for all } s \in \times_{i \in I} S_i\}.$$

Then,  $\psi$  satisfies the following properties:

- (1)  $\psi(\sigma, \tilde{\sigma}) = 1$ , if and only if,  $\sigma = \tilde{\sigma}$ .
- (2) If  $\text{supp}(\tilde{\sigma}) \subseteq \text{supp}(\sigma)$ , then  $\sigma(s) - \psi(\sigma, \tilde{\sigma}) \cdot \tilde{\sigma}(s) = 0$  for some  $s \in \text{supp}(\sigma)$ .
- (3) The function  $\psi(\cdot, \tilde{\sigma}) : \Delta(\times_{i \in I} S_i) \rightarrow \mathbb{R}$  is continuous, for all  $\tilde{\sigma} \in \Delta(\times_{i \in I} S_i)$ .

**Proof.** (1) Suppose that  $\psi(\sigma, \tilde{\sigma}) = 1$ . Then  $\sigma - 1 \cdot \tilde{\sigma} \geq 0$  and thus  $\sigma \geq \tilde{\sigma}$ . If  $\sigma(s') > \tilde{\sigma}(s')$  for some  $s' \in \times_{i \in I} S_i$ , then  $1 = \sum_{s \in \times_{i \in I} S_i} \sigma(s) > \sum_{s \in \times_{i \in I} S_i} \tilde{\sigma}(s) = 1$ , which is a contradiction. Therefore  $\sigma = \tilde{\sigma}$ . Conversely, suppose that  $\sigma = \tilde{\sigma}$ . Define  $\Psi_r := \{r \in \mathbb{R} : \sigma(s) - r \cdot \tilde{\sigma}(s) \geq 0, \text{ for all } s \in \times_{i \in I} S_i\}$ . Since  $\sigma - 1 \cdot \tilde{\sigma} = 0$ , then  $1 \in \Psi_r$ . Let  $\epsilon > 0$  and let  $s \in \text{supp}(\sigma) = \text{supp}(\tilde{\sigma})$ . Then  $\sigma(s) - (1 + \epsilon) \cdot \tilde{\sigma}(s) = -\epsilon \cdot \tilde{\sigma}(s) < 0$ . Hence,  $(1 + \epsilon) \notin \Psi_r$  for all  $\epsilon > 0$ . Therefore,  $\psi(\sigma, \tilde{\sigma}) = \sup_{r \in \mathbb{R}} \Psi_r = 1$ .

(2) Towards a contradiction, suppose that  $\text{supp}(\tilde{\sigma}) \subseteq \text{supp}(\sigma)$  and  $\sigma(s) - \psi(\sigma, \tilde{\sigma}) \cdot \tilde{\sigma}(s) > 0$  for all  $s \in \text{supp}(\sigma)$ . Let  $\bar{s} \in \arg \min\{\sigma(s) - \psi(\sigma, \tilde{\sigma}) \cdot \tilde{\sigma}(s) : s \in \text{supp}(\tilde{\sigma})\}$  and define  $r := \frac{(\sigma(\bar{s}) - \psi(\sigma, \tilde{\sigma}) \cdot \tilde{\sigma}(\bar{s}))}{\tilde{\sigma}(\bar{s})}$  and  $\psi'(\sigma, \tilde{\sigma}) := \psi(\sigma, \tilde{\sigma}) + r$ . Since  $\text{supp}(\tilde{\sigma})$  is finite,  $\bar{s}$  is well defined. Moreover, as  $\bar{s} \in \text{supp}(\tilde{\sigma}) \subseteq \text{supp}(\sigma)$ , then  $\sigma(\bar{s}) - \psi(\sigma, \tilde{\sigma}) \cdot \tilde{\sigma}(\bar{s}) > 0$ , hence  $r > 0$ , and thus  $\psi'(\sigma, \tilde{\sigma}) > \psi(\sigma, \tilde{\sigma})$ . Let  $s \in \times_{i \in I} S_i$ . If  $s \in (\times_{i \in I} S_i) \setminus \text{supp}(\sigma)$ , then  $\sigma(s) = \tilde{\sigma}(s) = 0$ , and thus  $\sigma(s) - \psi'(\sigma, \tilde{\sigma}) \cdot \tilde{\sigma}(s) = 0$ . If  $s \in \text{supp}(\sigma) \setminus \text{supp}(\tilde{\sigma})$ , then  $\sigma(s) > 0$  and  $\tilde{\sigma}(s) = 0$ , and thus  $\sigma(s) - \psi'(\sigma, \tilde{\sigma}) \cdot \tilde{\sigma}(s) > 0$ . If  $s \in \text{supp}(\tilde{\sigma})$ , then  $\sigma(s) - \psi(\sigma, \tilde{\sigma}) \cdot \tilde{\sigma}(s) \geq \sigma(\bar{s}) - \psi(\sigma, \tilde{\sigma}) \cdot \tilde{\sigma}(\bar{s}) > \sigma(\bar{s}) - \psi'(\sigma, \tilde{\sigma}) \cdot \tilde{\sigma}(\bar{s}) = \sigma(\bar{s}) - [\psi(\sigma, \tilde{\sigma}) + r] \cdot \tilde{\sigma}(\bar{s}) = \sigma(\bar{s}) - \psi(\sigma, \tilde{\sigma}) \cdot \tilde{\sigma}(\bar{s}) - r \cdot \tilde{\sigma}(\bar{s}) = 0$ . Consequently,  $\sigma(s) - \psi'(\sigma, \tilde{\sigma}) \cdot \tilde{\sigma}(s) \geq 0$  for all  $s \in \times_{i \in I} S_i$  and  $\psi'(\sigma, \tilde{\sigma}) > \psi(\sigma, \tilde{\sigma})$ , which contradicts the supremacy of  $\psi(\sigma, \tilde{\sigma})$ .

(3) Let  $\tilde{\sigma} \in \Delta(\times_{i \in I} S_i)$  and let  $(\sigma^k)_{k \in \mathbb{N}}$  be a sequence such that  $\lim_{k \rightarrow \infty} \sigma^k = \sigma$ . Then,  $\lim_{k \rightarrow \infty} \psi(\sigma^k, \tilde{\sigma}) = \psi(\lim_{k \rightarrow \infty} \sigma^k, \tilde{\sigma}) = \psi(\sigma, \tilde{\sigma})$ , and thus  $\psi(\cdot, \tilde{\sigma})$  is continuous. ■

**Lemma A.3.** Let  $(\sigma^k)_{k \in \mathbb{N}} \in (\Delta(\times_{i \in I} S_i))^{\mathbb{N}}$  be a sequence of mixed strategy profiles. Then, there exist a lexicographic probability measure  $\pi = (\pi^1, \dots, \pi^L) \in (\Delta(\times_{i \in I} S_i))^L$  and a sequence  $(r_n)_{n \in \mathbb{N}} = ((r_n^1, \dots, r_n^{L-1}))_{n \in \mathbb{N}} \in [(0, 1)^{L-1}]^{\mathbb{N}}$  with  $\lim_{n \rightarrow \infty} r_n = \vec{0}$  such that a subsequence  $(\sigma^{k_n})_{n \in \mathbb{N}}$  of  $(\sigma^k)_{k \in \mathbb{N}}$  satisfies  $\sigma^{k_n} = r_n \square \pi$  for all  $n \in \mathbb{N}$ .

**Proof.** Consider a subsequence  $(\sigma^{k_n})_{n \in \mathbb{N}}$  of  $(\sigma^k)_{k \in \mathbb{N}}$  that satisfies the following property: for every  $s \in \times_{i \in I} S_i$ , if  $\sigma^{k_n}(s) = 0$  for infinitely many indices  $n \in \mathbb{N}$ , then  $\sigma^{k_n}(s) = 0$  for all  $n \in \mathbb{N}$ . Then, there exists some index  $N \in \mathbb{N}$  such that the subsequence  $(\sigma^{k_n})_{n \geq N}$  of  $(\sigma^{k_n})_{n \in \mathbb{N}}$  satisfies the following property: for every  $s \in \times_{i \in I} S_i$ , if  $\sigma^{k_n}(s) = 0$ , then  $\sigma^{k_n}(s) = 0$  for all  $n \geq N$ . By the Bolzano–Weierstrass Theorem, there exists some convergent subsequence of  $(\sigma^{k_n})_{n \geq N}$ , denoted by  $(\sigma^k)_{k \in \mathbb{N}}$  for the sake of simplicity, with limit  $\pi^1 := \lim_{k \rightarrow \infty} \sigma^k$ .

Either  $\sigma^k = \pi^1$  infinitely often or  $\sigma^k = \pi^1$  finitely often. Suppose that  $\sigma^k = \pi^1$  infinitely often. Let  $(\sigma^{k_n})_{n \in \mathbb{N}}$  be a subsequence of  $(\sigma^k)_{k \in \mathbb{N}}$  such that  $\sigma^{k_n} = \pi^1$  for all  $n \in \mathbb{N}$ , let  $(r_n)_{n \in \mathbb{N}}$  be the empty sequence, and let  $\pi = (\pi^1)$ . It follows that  $\sigma^{k_n} = \pi^1 = r_n \square \pi$  for all  $n \in \mathbb{N}$ , which completes the proof in this case.

Otherwise, suppose that  $\sigma^k = \pi^1$  finitely often. Then, there exists  $N \in \mathbb{N}$  such that  $\sigma^k \neq \pi^1$  for all  $k \geq N$ . Let  $(\sigma^{k_n})_{n \in \mathbb{N}}$  be a subsequence of  $(\sigma^k)_{k \in \mathbb{N}}$  such that  $\sigma^{k_n} \neq \pi^1$  for all  $n \in \mathbb{N}$ . This subsequence is denoted by  $(\sigma^k)_{k \in \mathbb{N}}$  for the sake of simplicity. By Lemma A.2 (1),  $\psi(\sigma^k, \pi^1) \neq 1$  for all  $k \in \mathbb{N}$ . Consider the then well-defined sequence  $(\pi_k^2)_{k \in \mathbb{N}}$  given by

$$\pi_k^2 := \frac{\sigma^k - \psi(\sigma^k, \pi^1) \cdot \pi^1}{1 - \psi(\sigma^k, \pi^1)} \quad (3)$$

for all  $k \in \mathbb{N}$ . Note that for every  $s \in \bigtimes_{i \in I} S_i$  and for each  $k \in \mathbb{N}$ , if  $\sigma^k(s) = 0$ , then  $\pi^1(s) = 0$  and thus  $\pi_k^2(s) = 0$ . It follows that  $\text{supp}(\pi_k^2) \subseteq \text{supp}(\sigma^k)$  for all  $k \in \mathbb{N}$ . In addition, Lemma A.2 (2) ensures that for every  $k \in \mathbb{N}$ , there exists  $s \in \text{supp}(\sigma^k)$  such that  $\sigma^k(s) - \psi(\sigma^k, \pi^1) \cdot \pi^1(s) = 0$ , and thus  $s \notin \text{supp}(\pi_k^2)$ . Consequently,  $\text{supp}(\pi_k^2) \subsetneq \text{supp}(\sigma^k)$  for all  $k \in \mathbb{N}$ .

Eq. (3) can be rewritten as

$$\sigma^k = \psi(\sigma^k, \pi^1) \cdot \pi^1 + [1 - \psi(\sigma^k, \pi^1)] \cdot \pi_k^2 \quad (4)$$

for all  $k \in \mathbb{N}$ , where  $\psi(\sigma^k, \pi^1) \in (0, 1)$ . Lemma A.2 (3) and Lemma A.2 (1) ensure that  $\lim_{k \rightarrow \infty} \psi(\sigma^k, \pi^1) = \psi(\lim_{k \rightarrow \infty} \sigma^k, \pi^1) = \psi(\pi^1, \pi^1) = 1$ . Consider the sequence  $(r_k^1)_{k \in \mathbb{N}}$  defined by

$$r_k^1 := 1 - \psi(\sigma^k, \pi^1) \quad (5)$$

for all  $k \in \mathbb{N}$ , where  $\lim_{k \rightarrow \infty} r_k^1 = 1 - \lim_{k \rightarrow \infty} \psi(\sigma^k, \pi^1) = 0$ . Eqs. (4) and (5) imply that

$$\sigma^k = (1 - r_k^1) \cdot \pi^1 + r_k^1 \cdot \pi_k^2 \quad (6)$$

for all  $k \in \mathbb{N}$ .

By similar reasoning applied to the sequence  $(\pi_k^2)_{k \in \mathbb{N}}$ , it follows that there exists a convergent subsequence  $(\pi_{k_n}^2)_{n \in \mathbb{N}}$  of  $(\pi_k^2)_{k \in \mathbb{N}}$ , also denoted as  $(\pi_k^2)_{k \in \mathbb{N}}$  for the sake of simplicity, with limit  $\pi_2 := \lim_{k \rightarrow \infty} \pi_k^2$ . Either  $\pi_k^2 = \pi^2$  infinitely often or  $\pi_k^2 = \pi^2$  finitely often.

Suppose that  $\pi_k^2 = \pi^2$  infinitely often. Let  $(\pi_{k_n}^2)_{n \in \mathbb{N}}$  be a subsequence of  $(\pi_k^2)_{k \in \mathbb{N}}$  such that  $\pi_{k_n}^2 = \pi^2$  for all  $n \in \mathbb{N}$ , let  $(r_n)_{n \in \mathbb{N}} = (r_{k_n}^1)_{n \in \mathbb{N}}$  and let  $\pi = (\pi^1, \pi^2)$ . Eq. (6) ensures that

$$\sigma^{k_n} = (1 - r_{k_n}^1) \cdot \pi^1 + r_{k_n}^1 \cdot \pi^2 = r_n \square \pi$$

for all  $n \in \mathbb{N}$ , which completes the proof in this case.

Otherwise, suppose that  $\pi_k^2 = \pi^2$  finitely often. There exist sequences  $(\pi_k^3)_{k \in \mathbb{N}}$  and  $(r_k^2)_{k \in \mathbb{N}}$  such that the following properties hold:

$$\pi_k^2 = (1 - r_k^2) \cdot \pi^2 + r_k^2 \cdot \pi_k^3 \quad (7)$$

$$\pi_k^3 := \frac{\pi_k^2 - \psi(\pi_k^2, \pi^2) \cdot \pi^2}{1 - \psi(\pi_k^2, \pi^2)} \text{ and } \text{supp}(\pi_k^3) \subsetneq \text{supp}(\pi_k^2) \text{ for all } k \in \mathbb{N}$$

$$r_k^2 := 1 - \psi(\pi_k^2, \pi^2) \text{ for all } k \in \mathbb{N} \text{ and } \lim_{k \rightarrow \infty} r_k^2 = 0.$$

Eqs. (6) and (7) imply that

$$\sigma^k = (1 - r_k^1) \cdot \pi^1 + r_k^1 \cdot [(1 - r_k^2) \cdot \pi^2 + r_k^2 \cdot \pi_k^3]. \quad (8)$$

Iterating the same reasoning for the sequences  $(\pi_k^l)_{k \in \mathbb{N}}$  for  $l \geq 3$  guarantees that there exist a lexicographic level  $L \in \mathbb{N}$ ,  $\pi = (\pi^1, \dots, \pi^L) \in (\Delta(\bigtimes_{i \in I} S_i))^L$ , and  $(r_n)_{n \in \mathbb{N}} \in [(0, 1)^{L-1}]^{\mathbb{N}}$  such that  $\lim_{n \rightarrow \infty} r_n = \vec{0}$  and  $\sigma^{k_n} = r_n \square \pi$  for all  $n \in \mathbb{N}$ . Note that the iterative process necessarily terminates after finitely many rounds, since the set  $\bigtimes_{i \in I} S_i$  is finite and  $\text{supp}(\sigma^k) \supsetneq \text{supp}(\pi_k^2) \supsetneq \text{supp}(\pi_k^3) \supsetneq \dots$  for all  $k \in \mathbb{N}$ . ■

Equipped with Lemmas A.1, A.2, A.3, we can now proceed to formally establish Lemma 1.

**Proof ( $\Rightarrow$ ):** Suppose that  $\sigma$  constitutes a perfect equilibrium of  $\Gamma$ . Then, there exists a sequence of tuples of mixed strategies  $(\sigma^k)_{k \in \mathbb{N}}$  such that properties (i), (ii), and (iii) of Definition 6 hold. By Lemma A.3, there exists a lexicographic probability measure  $\pi = (\pi^1, \dots, \pi^L) \in (\Delta(\bigtimes_{i \in I} S_i))^L$  and a sequence  $(r_n)_{n \in \mathbb{N}} = ((r_n^1, \dots, r_n^{L-1}))_{n \in \mathbb{N}} \in [(0, 1)^{L-1}]^{\mathbb{N}}$  with  $\lim_{n \rightarrow \infty} r_n = \vec{0}$  such that some subsequence  $(\sigma^{k_n})_{n \in \mathbb{N}}$  of  $(\sigma^k)_{k \in \mathbb{N}}$  can

be expressed as  $\sigma^{k_n} = r_n \square \pi$  for all  $n \in \mathbb{N}$ . For every  $i \in I$ , define the lexicographic conjecture  $\beta_i := \text{marg}_{S_{-i}} \pi$ . We show that  $(\sigma^k)_{k \in \mathbb{N}}$ ,  $(\beta_i)_{i \in I}$ , and  $\pi$  satisfy properties (a), (b), (c), (d), and (e) of Definition 7.

First, note that property (e) of Definition 7 is directly satisfied. Since  $\sigma^{k_n} = r_n \square \pi$  is a product measure for all  $n \in \mathbb{N}$ , it follows that  $\pi$  is a tuple of product measures. Consequently,

$$\begin{aligned} \beta_i &= \text{marg}_{S_{-i}} \pi = \bigotimes_{j \in I \setminus \{i\}} \text{marg}_{S_j} \pi \\ &= \bigotimes_{j \in I \setminus \{i\}} \text{marg}_{S_j} \text{marg}_{S_{-i}} \pi = \bigotimes_{j \in I \setminus \{i\}} \text{marg}_{S_j} \beta_i \end{aligned}$$

for all  $i \in I$ , which yields property (d) of Definition 7. Moreover, property (i) ensures that  $\sigma = \lim_{n \rightarrow \infty} \sigma^{k_n} = \lim_{n \rightarrow \infty} (r_n \square \pi) = \pi^1$ . Hence,

$$\sigma_i = \text{marg}_{S_i} \sigma = \text{marg}_{S_i} \pi^1 = \text{marg}_{S_i} \text{marg}_{S_{-i}} \pi^1 = \text{marg}_{S_i} \beta_i^1$$

for all  $i \in I$  and all  $j \in I \setminus \{i\}$ , which establishes property (b) of Definition 7. Furthermore, property (ii) guarantees that  $\sigma_i^{k_n}$  has full support for all  $i \in I$  and for all  $n \in \mathbb{N}$ . Thus,  $\pi$  and hence  $\beta_i$  is cautious for all  $i \in I$ , which establishes property (a) of Definition 7. Finally, let  $s_i \in \text{supp}(\sigma_i)$ . By property (iii),  $s_i$  is a best response to  $\sigma_{-i}^{k_n} = \text{marg}_{S_{-i}} (r_n \square \pi) = r_n \square \beta_i$  for all  $n \in \mathbb{N}$ . By Lemma A.1 (ii),  $s_i$  is lex-optimal given  $\beta_i$ , which corresponds to property (c) of Definition 7. Therefore,  $\sigma = (\sigma_i)_{i \in I}$  constitutes a lexicographic perfect equilibrium of  $\Gamma$ .

( $\Leftarrow$ ): Suppose that  $\sigma$  constitutes a lexicographic perfect equilibrium of  $\Gamma$ . Then, there exists a tuple of lexicographic conjectures  $\beta = (\beta_i)_{i \in I}$  and a lexicographic product measure  $\pi = (\pi^1, \dots, \pi^L)$  satisfying properties (a), (b), (c), (d), and (e) of Definition 7. Consider the sequence  $(r_n)_{n \in \mathbb{N}} = ((\frac{1}{n+1}, \dots, \frac{1}{n+1}))_{n \in \mathbb{N}} \in [(0, 1)^{L-1}]^{\mathbb{N}}$ . Note that  $\lim_{n \rightarrow \infty} r_n = \vec{0}$ . For every  $i \in I$  and for every  $n \in \mathbb{N}$ , define  $\sigma_i^n := \text{marg}_{S_i} (r_n \square \pi)$  and  $\sigma^n := (\sigma_i^n)_{i \in I}$ . We show that there exists a subsequence of  $(\sigma^n)_{n \in \mathbb{N}}$  satisfying properties (i), (ii), (iii) of Definition 6.

Let  $i \in I$  be some player. Since  $r_n \square \pi$  is a product measure and properties (b) and (e) hold,

$$\begin{aligned} \lim_{n \rightarrow \infty} \sigma_i^n &= \lim_{n \rightarrow \infty} \text{marg}_{S_i} (r_n \square \pi) = \lim_{n \rightarrow \infty} \text{marg}_{S_i} \text{marg}_{S_{-i}} (r_n \square \pi) \\ &= \lim_{n \rightarrow \infty} \text{marg}_{S_i} (r_n \square \text{marg}_{S_{-i}} \pi) = \lim_{n \rightarrow \infty} \text{marg}_{S_i} (r_n \square \beta_i) \\ &= \text{marg}_{S_i} \lim_{n \rightarrow \infty} (r_n \square \beta_i) = \text{marg}_{S_i} \beta_i^1 = \sigma_i \end{aligned}$$

for all  $j \in I$  such that  $i \neq j$ . This establishes property (i) of Definition 6. In addition, let  $j \in I \setminus \{i\}$ ,  $s_j \in S_j$ , and  $n \in \mathbb{N}$ . List (a) ensures that there exists a level  $l^* \in \{1, \dots, L\}$  such that  $\text{marg}_{S_j} \beta_i^{l^*}(s_j) > 0$ . It follows that

$$\begin{aligned} \sigma_j^n(s_j) &= \text{marg}_{S_j} (r_n \square \pi)(s_j) = \text{marg}_{S_j} \text{marg}_{S_{-j}} (r_n \square \pi)(s_j) \\ &= \text{marg}_{S_j} (r_n \square \text{marg}_{S_{-j}} \pi)(s_j) = \text{marg}_{S_j} (r_n \square \beta_j)(s_j) > 0. \end{aligned}$$

Hence,  $\text{supp}(\sigma_j^n) = S_j$ , which yields property (ii) of Definition 6. Besides, let  $s_i \in \text{supp}(\sigma_i)$ . List (c) ensures that  $s_i$  is lex-optimal given  $\beta_i$ . By Lemma A.1 (iii), there exists some subsequence  $(r_{n_k})_{k \in \mathbb{N}}$  of  $(r_n)_{n \in \mathbb{N}}$  and some index  $K \in \mathbb{N}$  such that  $u_i(s_i, r_{n_k} \square \beta_i) \geq u_i(s_i', r_{n_k} \square \beta_i)$  for all  $k \geq K$  and for all  $s_i' \in S_i$ . List (e) guarantees that

$$\begin{aligned} r_{n_k} \square \beta_i &= (r_{n_k} \square \text{marg}_{S_{-i}} \pi) = \text{marg}_{S_{-i}} (r_{n_k} \square \pi) \\ &= \bigotimes_{j \in I \setminus \{i\}} \text{marg}_{S_j} (r_{n_k} \square \pi) = \bigotimes_{j \in I \setminus \{i\}} \sigma_j^{n_k}. \end{aligned}$$

Hence,  $s_i$  is a best response to  $\sigma_{-i}^{n_k}$  for all  $k \geq K$ , i.e. the subsequence  $(\sigma_{-i}^{n_k})_{k \geq K}$  satisfies property (iii) of Definition 6. Consequently, the subsequence  $(\sigma^{n_k})_{k \geq K}$  satisfies properties (i), (ii), (iii) of Definition 6. Therefore,  $\sigma$  constitutes a perfect equilibrium of  $\Gamma$ . ■

## References

- Asheim, G.B., 2001. Proper rationalizability in lexicographic beliefs. *Internat. J. Game Theory* 30, 453–478.
- Asheim, G.B., 2002. On the epistemic foundation for backward induction. *Math. Social Sci.* 44, 121–144.

- Asheim, G.B., Perea, A., 2005. Sequential and quasi-perfect rationalizability in extensive games. *Games Econ. Behav.* 53, 15–42.
- Aumann, R.J., 1974. Subjectivity and correlation in randomized strategies. *J. Math. Econom.* 1, 67–96.
- Aumann, R.J., 1976. Agreeing to disagree. *Ann. Statist.* 4, 1236–1239.
- Aumann, R.J., Brandenburger, A., 1995. Epistemic conditions for Nash equilibrium. *Econometrica* 63, 1161–1180.
- Bach, C.W., Cabessa, J., 2017. Limit-agreeing to disagree. *J. Log. Comput.* 27, 1169–1187.
- Bach, C.W., Perea, A., 2013. Agreeing to disagree with lexicographic prior beliefs. *Math. Social Sci.* 66, 129–133.
- Bach, C.W., Perea, A., 2020. Generalized Nash equilibrium without common belief in rationality. *Econom. Lett.* 186, 1–6.
- Bach, C.W., Tsakas, E., 2014. Pairwise epistemic conditions for Nash equilibrium. *Games Econ. Behav.* 85, 48–59.
- Barelli, P., 2009. Consistency of beliefs and epistemic conditions for Nash and correlated equilibria. *Games Econ. Behav.* 67, 363–375.
- Battigalli, P., Siniscalchi, M., 1999. Hierarchies of conditional beliefs and interactive epistemology in dynamic games. *J. Econom. Theory* 88, 188–230.
- Battigalli, P., Siniscalchi, M., 2002. Strong belief and forward induction reasoning. *J. Econom. Theory* 106, 356–391.
- Blume, L., Brandenburger, A., Dekel, E., 1991a. Lexicographic probabilities and choice under uncertainty. *Econometrica* 59, 61–79.
- Blume, L., Brandenburger, A., Dekel, E., 1991b. Lexicographic probabilities and equilibrium refinements. *Econometrica* 59, 81–98.
- Bonanno, G., 2018. Behavior and deliberation in perfect-information games: Nash equilibrium and backward induction. *Internat. J. Game Theory* 47, 1001–1032.
- Bonanno, G., Nehring, K., 1997. Agreeing to Disagree: A Survey. Mimeo.
- Börgers, T., 1994. Weak dominance and approximate common knowledge. *J. Econom. Theory* 64, 265–276.
- Brandenburger, A., 1992a. Knowledge and equilibrium in games. *J. Econ. Perspect.* 6, 83–101.
- Brandenburger, A., 1992b. Lexicographic probabilities and iterated admissibility. In: Dasgupta, P., et al. (Eds.), *Economic Analysis of Markets and Games*. MIT Press, pp. 282–290.
- Brandenburger, A., Dekel, E., 1987. Common knowledge with probability 1. *J. Math. Econom.* 16, 237–245.
- Brandenburger, A., Friedenberg, A., Keisler, J., 2008. Admissibility in games. *Econometrica* 76, 307–352.
- Catonini, E., De Vito, N., 2018. Cautious Belief and Iterated Admissibility. Mimeo.
- Catonini, E., De Vito, N., 2020. Weak belief and permissibility. *Games Econ. Behav.* 120, 154–179.
- Chen, Y.-C., Lehrer, E., Li, J., Samet, D., Shmaya, E., 2015. Agreeing to disagree and dutch books. *Games Econ. Behav.* 93, 108–116.
- Contreras-Tejada, P., Scarpa, G., Kubicki, A.M., Brandenburger, A., La Mura, P., 2021. Observers of quantum systems cannot agree to disagree. *Nature Commun.* 12, article number: 7021.
- van Damme, E., 1984. A relationship between perfect equilibria in extensive form games and proper equilibria in normal form games. *Int. J. Game Theory* 13, 1–13.
- Dégremont, C., Roy, O., 2012. Agreement theorems in dynamic-epistemic logic. *J. Physiol (London)* 41, 735–764.
- Dekel, E., Friedenberg, A., Siniscalchi, M., 2016. Lexicographic beliefs and assumption. *J. Econom. Theory* 163, 955–985.
- Demey, L., 2014. Agreeing to disagree in probabilistic dynamic epistemic logic. *Synthese* 191, 409–438.
- Dominiak, A., Lefort, J.-P., 2015. Agreeing to disagree type results under ambiguity. *J. Math. Econom.* 61, 119–129.
- Gale, D., 1953. A theory of  $n$ -person games with perfect information. *Proc. Natl. Acad. Sci.* 39, 496–501.
- Gizatulina, A., Hellman, Z., 2019. No trade and yes trade theorems for heterogeneous priors. *J. Econom. Theory* 182, 161–184.
- Govindan, S., Klumpp, T., 2003. Perfect equilibrium and lexicographic beliefs. *Int. J. Game Theory* 31, 229–243.
- Halpern, J., 2010. Lexicographic probability, conditional probability, and nonstandard probability. *Games Econ. Behav.* 68, 155–179.
- Hammond, P., 1994. Elementary non-Archimedean representation of probability for decision theory and games. In: Humphrey, P. (Ed.), *Patrick Suppes: Scientific Philosopher*, 1: Probability and Probabilistic Causality. Kluwer, pp. 25–61.
- Harsanyi, J.C., 1967–68. Games of incomplete information played by Bayesian players. Part I, II, III. *Manage. Sci.* 14, 159–182, 320–334, 486–502.
- Harsanyi, J.C., 1973. Games with randomly disturbed payoffs: A new rationale for mixed-strategy equilibrium points. *Internat. J. Game Theory* 2, 1–23.
- Harsanyi, J.C., Selten, R., 1988. *A General Theory of Equilibrium Selection in Games*. MIT Press.
- Heifetz, A., Meier, M., Schipper, B.C., 2013. Unawareness, beliefs and speculative trade. *Games Econ. Behav.* 77, 100–121.
- Hellman, Z., 2013. Almost common priors. *Internat. J. Game Theory* 42, 399–410.
- Hellman, Z., Samet, D., 2012. How common are common priors? *Games Econ. Behav.* 74, 517–525.
- Kreps, D., Ramey, G., 1987. Structural consistency, consistency, and sequential rationality. *Econometrica* 55, 1331–1348.
- Kreps, D., Wilson, R., 1982. Sequential equilibria. *Econometrica* 50, 863–894.
- Lee, B.S., 2016. Admissibility and assumption. *J. Econom. Theory* 163, 42–72.
- Lehrer, E., Samet, D., 2014. Belief consistency and trade consistency. *Games Econ. Behav.* 83, 165–177.
- Liu, Y., 2019. Two tales of epistemic models. *Thought J. Philos.* 8, 291–302.
- Mailath, G.J., Samuelson, L., Swinkels, J.M., 1997. How proper is sequential equilibrium? *Games Econ. Behav.* 18, 193–218.
- Ménager, L., 2012. Agreeing to Disagree: A Review. Mimeo.
- Milgrom, P., Stokey, N., 1982. Information, trade and common knowledge. *J. Econom. Theory* 26, 17–27.
- Myerson, R.B., 1978. Refinements of the Nash equilibrium concept. *Internat. J. Game Theory* 7, 73–80.
- Nash, J., 1950. Equilibrium points in  $N$ -person games. *Proc. Natl. Acad. Sci.* 36, 48–49.
- Nash, J., 1951. Non-cooperative games. *Ann. Mat.* 54, 286–295.
- Pacuit, E., 2018. Agreement Theorems with Qualitative Conditional Probability. Mimeo.
- Perea, A., 2007. A one-person doxastic characterization of Nash strategies. *Synthese* 158, 1251–1271.
- Perea, A., 2012. *Epistemic Game Theory: Reasoning and Choice*. Cambridge University Press.
- Rényi, A., 1995. On a new axiomatic theory of probability. *Acta Math. Acad. Sci. Hung.* 6, 285–335.
- Samuelson, L., 1992. Dominated strategies and common knowledge. *Games Econ. Behav.* 4, 284–313.
- Selten, R., 1975. Reexamination of the perfectness concept for equilibrium points in extensive games. *Internat. J. Game Theory* 4, 25–55.
- Stahl, D.O., 1995. Lexicographic rationalizability and iterated admissibility. *Econom. Lett.* 47, 155–159.
- Stuart, H.W., 1997. Common belief of rationality in the finitely repeated prisoners' dilemma. *Games Econ. Behav.* 19, 133–143.
- Tarbusch, B., 2016. Counterfactuals in agreeing to disagree type results. *Math. Social Sci.* 84, 125–133.
- Tsakas, E., 2014. Epistemic equivalence of extended belief hierarchies. *Games Econ. Behav.* 86, 126–144.
- Tsakas, E., 2018. Agreeing to disagree with conditional probability systems. *B. E. J. Theor. Econ.* 18, 1–7.
- Yang, C., 2015. Weak assumption and iterated admissibility. *J. Econom. Theory* 158, 87–101.