

# Valuing others' opinions: Preference, belief and reliability dynamics

Sujata Ghosh<sup>1</sup> and Katsuhiko Sano<sup>2</sup>

<sup>1</sup> Indian Statistical Institute, Chennai, India

<sup>2</sup> Department of Philosophy, Graduate School of Letters, Hokkaido University, Sapporo, Japan  
sujata@isichennai.res.in, katsuhiko.sano@gmail.com

Keywords: Preference, Belief, Reliability, Hybrid Logic, Public Announcement Logic, Propositional Dynamic Logic

Abstract: Deliberation often leads to changes in preferences and beliefs of an agent, influenced by the opinions of others, depending on how reliable these agents are according to the agent under consideration. Sometimes, it also leads to changes in the opposite direction, that is, reliability over agents gets updated depending on their preferences and/or beliefs. There are various formal studies of preference and belief change based on reliability and/or trust, but not the other way around – this work contributes to the formal study of the latter aspect, that is, on reliability change based on agent preferences. In process, some policies of preference change based on agent reliabilities are also discussed. A two-dimensional hybrid language is proposed to describe such processes, and axiomatisations and decidability are discussed.

## 1 INTRODUCTION

Deliberation forms an important component in any decision-making process. It is basically a conversation through which individuals provide their opinions regarding certain issues, give preferences among possible choices, justify these preferences. This process may lead to changes in their opinions, because they are influenced by one another. A factor that sometimes plays a big role in enforcing such changes, is the amount of reliability the agents have on one another's opinions. Such reliabilities may change as well through this process of deliberation, e.g. on hearing someone else's preferences about a certain issue, one can start or stop relying on that person's opinion. One may tend to unfriend certain friends hearing about their preferences regarding certain issues (e.g. Helen De Cruz's recent remarks in her article 'Being Friends with a Brexiter?' in the Philosophers On series of the Daily Nous blog<sup>1</sup>).

Formal studies on preferences (cf. (Arrow et al., 2002; Endriss, 2011)) and trust (cf. (Liau, 2003; Demolombe, 2004; Herzig et al., 2010)) abound in the literature on logic in artificial intelligence. Recently, there has been work on relating the notions of belief and trust, e.g. about agents changing their beliefs based on another agent's announcement depending on how trustworthy that agent is about the issue

in question (e.g. see (Lorini et al., 2014)). And, also on relating preference and reliability, e.g. about agents changing their preferences based on another agent's preferences, on whom he or she relies the most (Ghosh and Velázquez-Quesada, 2015a; Ghosh and Velázquez-Quesada, 2015b). A pertinent issue that arises in this context is: an agent's assessment of another individual's reliability might change as well. How would one model that? This work precisely provides a way to answer this question. We focus on reliability changes based on (public) announcement of individual preferences and we provide formal frameworks to describe such changes. In process, we also provide some policies of preference change as well. Note that the notion of reliability considered here is not topic-based (in contrast to the notion of trust described in (Lorini et al., 2014)) but deals with only comparative judgements about agents (cf. Section 2 for details). The following provides an apt example of the situations we would like to model:

*Our Running Example:* Consider three flat-mates Isabella, John and Ken discussing about redecorating their house and they were wondering whether to put a print of Monet's picture on the left wall or on the right wall of the living room. Isabella and Ken prefer to put it on the right wall, while John wants to put it on the left. Isabella has more faith in John's taste than on hers and Ken's, and John has more faith in Isabella's taste than on his' and Ken's.

<sup>1</sup><http://dailynous.com/2016/06/28/philosophers-on-brexiter/#DeCruz>

Ken has full faith in his own taste. As long as Isabella’s and John’s preferences are different and each think that the other’s taste is better (by taking their preferences into consideration), the three flatmates would never reach an agreement. But it so happens that on hearing about John’s and Ken’s choices, Isabella starts relying more on Ken, whereas even after hearing about Isabella’s and Ken’s choices John’s reliability attribution to Isabella does not change.

To model such situations we introduce a two-dimensional hybrid logic framework extending the basic logic proposed in (Ghosh and Velázquez-Quesada, 2015a; Ghosh and Velázquez-Quesada, 2015b) in the line of those developed in (Gargov et al., 1987; Sano, 2010; Seligman et al., 2013). We add dynamic operators to model preference and reliability changes. The main novelty of this work is that reliability changing policies based on agent preferences are introduced and studied formally, which has not been dealt with before. In addition, reliabilities are modelled (more naturally) as total pre-orders (instead of total orders (Ghosh and Velázquez-Quesada, 2015a; Ghosh and Velázquez-Quesada, 2015b)), and preference changing policies are modified accordingly. The proposed logic is expressive enough to deal with both these kinds of changes.

## 2 TWO-DIMENSIONAL HYBRID LOGIC

Let us first motivate our assumptions on preference and reliability orders that we make below in the lines of (Ghosh and Velázquez-Quesada, 2015a). As mentioned earlier, we are modelling situations akin to joint deliberation where agents announce their preferences. Each agent can change her preferences upon getting information about the other agents’ preferences, influenced by her reliability over agents (including herself, so she might consider herself as more reliable than some agents but also as less reliable than some others). Agents can also change their opinions regarding how reliable they think the other agents are in comparison to themselves, influenced by the announced preferences of those agents.

The agents’ preferences are represented by binary relations (as in (Arrow et al., 2002; Grüne-Yanoff and Hansson, 2009) and further references therein), which is typically assumed to be reflexive and transitive. This paper also follows this ordinary assumption, and so, we note that this assumption do allow the possibility of incomparable worlds.

The notion of *reliability* is related to that of *trust*, a well-studied concept (e.g., (Falcone et al., 2008)), with several proposals for its formal representation, e.g. an attitude of an agent who believes that another agent has a given property (Falcone and Castelfranchi, 2001). One also says that “an agent  $i$  trusts agent  $j$ ’s judgement about  $\varphi$ ” (called “trust on credibility” in (Demolombe, 2001)). Trust can also be defined in terms of other attitudes, such as knowledge, beliefs, intentions and goals (e.g., (Demolombe, 2001; Herzig et al., 2010)), or as a semantic primitive, typically by means of a *neighbourhood function* (Liau, 2003). Some others (e.g., (Lorini et al., 2014)) deal with *graded trust*.

Reliability as discussed here is closer to the notion of trust in (Holliday, 2010), where it is understood as an ordering among sets of sources of information (cf. the discussion in (Goldman, 2001)). Such a notion of reliability does not yield any *absolute* judgements (“ $i$  relies on  $j$ ’s judgement [about  $\varphi$ ]”), but only *comparative* ones (“for  $i$ , agent  $j$ ’ is at least as reliable as agent  $j$ ”). For the purposes of this work, similar to (Ghosh and Velázquez-Quesada, 2015a), such comparative judgements suffice.

In contrast to (Ghosh and Velázquez-Quesada, 2015a), our reliability relation is assumed to be a reflexive, transitive and total relation. Reflexivity and transitivity are, more often than not, natural requirements for an ordering and totality disallows incomparability, as before. The changes in reliability for an agent depend on the information assimilated (similar to approaches like (Rodenhäuser, 2014)), in particular, about the other agents’ preferences.

The focus of this work is joint deliberation, so let  $A$  be a *finite non-empty* set of agents ( $|A| = n \geq 2$ ).

**Definition 1.** A *PR (preference/reliability) frame*  $F$  is a tuple  $(W, \{\leq_i, \leq_i\}_{i \in A})$  where **(1)**  $W$  is a finite non-empty set of worlds; **(2)**  $\leq_i \subseteq W \times W$  is a preorder (i.e., a reflexive and transitive relation), agent  $i$ ’s *preference relation* among worlds in  $W$  ( $u \leq_i v$  is read as “world  $v$  is at least as preferable as world  $u$  for agent  $i$ ”); **(3)**  $\leq_i \subseteq A \times A$  is a total pre-order (i.e., a connected pre-order), agent  $i$ ’s *reliability relation* among agents in  $A$  ( $j \leq_i k$  is read as “agent  $k$  is at least as reliable as agent  $j$  for agent  $i$ ”). Let  $mr(i)$  denote the set of all maximally reliable agents for  $i$ .

We define  $u <_i v$  (“ $u$  is less preferred than  $v$  for agent  $i$ ”) as  $u \leq_i v$  and  $v \not\leq_i u$ , and  $u \approx_i v$  (“ $u$  and  $v$  are equally preferred for agent  $i$ ”) as  $u \leq_i v$  and  $v \leq_i u$ . Moreover,  $j <_i k$  (“ $j$  is less reliable than  $k$  for agent  $i$ ”) is defined as  $j \leq_i k$  and  $k \not\leq_i j$ , and  $j \approx_i k$  (“ $j$  and  $k$  are equally reliable for agent  $i$ ”) as  $j \leq_i k$  and  $k \leq_i j$ .

**Example 1.** Recall the example in Section 1. Put  $A = \{i, j, k\}$ , where  $i$ ,  $j$ , and  $k$  represent Isabella, John,

and Ken, respectively. By denoting with  $w_x$  the world where ‘Monet’s picture is at wall  $x$ ’ ( $x = l, r$ ), the example’s situation can be represented by a *PR* frame  $F_{exp} = (\{w_l, w_r\}, \{\leq_y, \leq_x\}_{y \in A})$  in which the preference orders are given by:  $w_l <_i w_r$ ,  $w_r <_j w_l$  and  $w_l <_k w_r$ , and the reliability orders are given by:  $i \approx_i k <_i j$ ,  $j \approx_j k <_j i$  and  $j \approx_k i <_k k$ .

In (Ghosh and Velázquez-Quesada, 2015a), Ghosh and Velázquez-Quesada propose a language to talk about the preference changes and their effects. Following the semantic idea of (Seligman et al., 2013), we extend their syntax for the static language into a two-dimensional syntax with the help of dependent product of two hybrid logics (Sano, 2010). Let  $P$  be a countable infinite set of propositional variables,  $N_1 = \{\mathbf{a}, \mathbf{b}, \mathbf{c}, \dots\}$  be a countable infinite set of world-nominals (syntactic names for worlds) and let  $N_2 = \{\mathbf{i}, \mathbf{j}, \mathbf{k}, \dots\}$  be a countable infinite set of agent-nominals (syntactic names for agents).

**Definition 2** (Language  $\mathcal{HL}$ ). Formulas  $\varphi, \psi, \dots$  (read  $\varphi$  as “the current agent satisfies the property  $\varphi$  in the current state” or indexically as “I am  $\varphi$  in the current state”) and relational expressions (or program terms)  $\pi, \rho, \dots$  of the language  $\mathcal{HL}$  are given by

$$\begin{aligned} \varphi, \psi &::= p \mid \mathbf{a} \mid \mathbf{i} \mid \neg \varphi \mid \varphi \vee \psi \mid @_i \varphi \mid @_{\mathbf{a}} \varphi \mid \langle \pi \rangle \varphi, \\ \pi, \rho &::= 1_W \mid \leq \mid \geq \mid 1_A \mid \sqsubseteq_{\mathbf{k}} \mid \sqsupseteq_{\mathbf{k}} \mid -\alpha \\ &\quad \pi \cup \rho \mid \pi \cap \rho \mid (\pi, \mathbf{j}) \sqcap_i (\rho, \mathbf{k}) \mid ?(\varphi, \psi), \end{aligned}$$

where  $p \in P$ ,  $\mathbf{a} \in N_1$ ,  $\mathbf{i}, \mathbf{j}, \mathbf{k} \in N_2$ ,  $\alpha \in \{1_W, \leq, \geq\} \cup \{1_A, \sqsubseteq_{\mathbf{k}}, \sqsupseteq_{\mathbf{k}} \mid \mathbf{k} \in N_2\}$ . Propositional constants ( $\top, \perp$ ), other Boolean connectives ( $\wedge, \rightarrow, \leftrightarrow$ ) and the dual modal universal operators  $[\pi]$  are defined as usual, e.g.  $[\pi]\varphi := \neg \langle \pi \rangle \neg \varphi$ . Moreover, we define  $\langle \leq \rangle \varphi$  as  $\langle \leq \cap - \geq \rangle \varphi$  and  $\langle \sqsubseteq_{\mathbf{k}} \rangle \varphi$  as  $\langle \sqsubseteq_{\mathbf{k}} \cap - \sqsupseteq_{\mathbf{k}} \rangle \varphi$ , respectively. Finally,  $?\varphi$  is defined as the program term  $?( \varphi, \varphi )$ .

We note that  $@_{\mathbf{a}}\varphi$  is read as “the current agent satisfies  $\varphi$  in the world named by  $\mathbf{a}$ ” and  $@_i\varphi$  as “agent  $\mathbf{i}$  satisfies  $\varphi$  in the current world.” The set of relational expressions contains the constants  $1_W, 1_A$  (the global relations, whose corresponding operators mean “for all states” and “for all agents”, respectively), the preference and reliability relations ( $\leq, \sqsubseteq_{\mathbf{k}}$ ), their respective converse relations ( $\geq, \sqsupseteq_{\mathbf{k}}$ ; cf. (Burgess, 1984; Goldblatt, 1992)), all the complements of the atomic relations, and an additional construct of the forms  $(\pi, \mathbf{j}) \sqcap_i (\pi', \mathbf{k})$  (needed for defining distributed preference later in Section 3, explained below) and  $?( \varphi, \psi )$  (a generalization of the test operator in (Harel et al., 2000), also explained below), and it is closed under union and intersection operations over relations.

The formulas are interpreted in terms of *world-agent* pairs below, and we may read  $\langle \leq \rangle \varphi$  as “in all states which the current agent considers as least as

good as the current state, the current agent satisfies  $\varphi$ ”. Moreover, we may read  $\langle \sqsubseteq_{\mathbf{k}} \rangle \varphi$  as “there is a more or equally reliable agent  $j$  than the current agent such that  $j$  satisfies  $\varphi$ , from agent  $\mathbf{k}$ ’s perspective.” For example,  $@_i \langle \sqsubseteq_{\mathbf{k}} \rangle \mathbf{j}$  can be read as “agent  $\mathbf{j}$  is more or equally reliable than agent  $\mathbf{i}$  from agent  $\mathbf{k}$ ’s perspective.”  $\langle \sqsupseteq_{\mathbf{k}} \rangle \varphi$  is read as “there is a less or equally reliable agent  $j$  than the current agent such that  $j$  satisfies  $\varphi$ , from agent  $\mathbf{k}$ ’s perspective.”

We note that the program construction  $?( \varphi, \psi )$  (check if the first element of a given pair of states satisfies  $\varphi$  and if the second does  $\psi$ ) is a generalization of the test operator in the standard (regular) propositional dynamic logic (Harel et al., 2000). So  $?\varphi := ?(\varphi, \varphi)$  enables us to check if both elements of a given pair satisfies  $\varphi$ . Moreover, the program construction  $(\pi, \mathbf{j}) \sqcap_i (\pi', \mathbf{k})$  enables us to define, as agent  $\mathbf{i}$ ’s relation between states, the distributed preference between agents  $\mathbf{j}$  and  $\mathbf{k}$ , i.e., the intersection of  $\mathbf{j}$ ’s preference and  $\mathbf{k}$ ’s preference. Together with the other program constructions, it is useful for providing the axiom system for the preference and reliability changing operations to be introduced in Section 3. The following two definitions establish what a model is and how formulas of  $\mathcal{HL}$  are interpreted over them.

**Definition 3** (PR model). A *PR* model is a tuple  $M = (F, V)$  where  $F = (W, \{\leq_i, \leq_x\}_{i \in A})$  is a *PR*-frame and  $V$  is a valuation function from  $P \cup N_1 \cup N_2$  to  $\mathcal{P}(W \times A)$  assigning a subset of the form  $\{w\} \times A$  to a world-nominal  $\mathbf{a} \in N_1$  and a subset of the form  $W \times \{i\}$  to an agent-nominal  $\mathbf{i} \in N_2$ . Throughout the paper, we denote the unique element in the first coordinate of  $V(\mathbf{a}) = \{w\} \times A$  and the second coordinate of  $V(\mathbf{i}) = W \times \{a\}$  by  $\underline{\mathbf{a}}$  and  $\underline{\mathbf{i}}$ , respectively.

**Definition 4** (Truth definition). Given a *PR*-model  $M$ , a satisfaction relation  $M, (w, i) \Vdash \varphi$ , and relations  $R_{\pi} \subseteq (W \times A)^2$  are defined by simultaneous induction by:

$$\begin{aligned} M, (w, i) \Vdash p &\text{ iff } (w, i) \in V(p), \\ M, (w, i) \Vdash \mathbf{a} &\text{ iff } w = \underline{\mathbf{a}}, \\ M, (w, i) \Vdash \mathbf{i} &\text{ iff } i = \underline{\mathbf{i}}, \\ M, (w, i) \Vdash \neg \varphi &\text{ iff } M, (w, i) \not\Vdash \varphi, \\ M, (w, i) \Vdash \varphi \vee \psi &\text{ iff } M, (w, i) \Vdash \varphi \text{ or } M, (w, i) \Vdash \psi, \\ M, (w, i) \Vdash @_{\mathbf{a}} \varphi &\text{ iff } M, (\underline{\mathbf{a}}, i) \Vdash \varphi, \\ M, (w, i) \Vdash @_i \varphi &\text{ iff } M, (w, \underline{\mathbf{i}}) \Vdash \varphi, \\ M, (w, i) \Vdash \langle \pi \rangle \psi &\text{ iff For some } (v, j) \in W \times A, \\ &\quad (w, i) R_{\pi} (v, j) \text{ and } M, (v, j) \Vdash \psi, \\ (w, i) R_{1_W} (v, j) &\text{ iff } w, v \in W \text{ and } i = j, \\ (w, i) R_{\leq} (v, j) &\text{ iff } w \leq_i v \text{ and } i = j, \\ (w, i) R_{\geq} (v, j) &\text{ iff } v \leq_i w \text{ and } i = j, \end{aligned}$$

$$\begin{aligned}
(w, i)R_{-\alpha}(v, j) & \text{ iff } ((w, i), (v, i)) \notin R_\alpha \text{ and } i = j \\
& \quad (\alpha \in \{1_W, \leq, \geq\}), \\
(w, i)R_{1_A}(v, j) & \text{ iff } w = v \text{ and } i, j \in A, \\
(w, i)R_{\sqsubseteq_{\mathbf{k}}}(v, j) & \text{ iff } w = v \text{ and } i \leq_{\mathbf{k}} j, \\
(w, i)R_{\supseteq_{\mathbf{k}}}(v, j) & \text{ iff } w = v \text{ and } j \leq_{\mathbf{k}} i, \\
(w, i)R_{-\beta}(v, j) & \text{ iff } w = v \text{ and } ((w, i), (w, j)) \notin R_\beta \\
& \quad (\beta \in \{1_A, \sqsubseteq_{\mathbf{k}}, \supseteq_{\mathbf{k}} \mid \mathbf{k} \in \mathbb{N}_2\}), \\
(w, i)R_{\pi \cup \rho}(v, j) & \text{ iff } (w, i)R_\pi(v, j) \text{ or } (w, i)R_\rho(v, j), \\
(w, i)R_{\pi \cap \rho}(v, j) & \text{ iff } (w, i)R_\pi(v, j) \text{ and } (w, i)R_\rho(v, j), \\
(w, i)R_{(\pi, \mathbf{j}) \sqcap (\rho, \mathbf{k})}(v, j) & \text{ iff } i = j = \mathbf{i} \text{ and } (w, \mathbf{j})R_\pi(v, \mathbf{j}) \\
& \quad \text{and } (w, \mathbf{k})R_\rho(v, \mathbf{k}) \\
(w, i)R_{\triangleright(\varphi, \psi)}(v, j) & \text{ iff } M, (w, i) \Vdash \varphi \text{ and } M, (v, j) \Vdash \psi.
\end{aligned}$$

We say that  $\varphi$  is *valid* in a *PR*-model  $M$  (written:  $M \Vdash \varphi$ ) if  $M, (w, i) \Vdash \varphi$  for all pairs  $(w, i)$  in  $M$ .

The logic  $\mathcal{HL}$  is so expressive that we can formalize the notion of belief as well as our preference and reliability dynamics introduced in the later sections. For example, following the idea found in (Boutilier, 1994), we can define the conditional belief operator  $B(\psi, \varphi)$  (read “under the condition that the current agent satisfies  $\psi$ , the current agent believes that she satisfies  $\varphi$ ” or “the current agent desires (or has a goal) that she satisfies  $\varphi$  under the condition that she satisfies  $\psi$ ”) by

$$B^\psi \varphi := [1_W]((\psi \wedge \varphi) \rightarrow \langle \leq \rangle (\psi \wedge \varphi \wedge [\leq](\psi \rightarrow \varphi))).$$

Then the unconditional belief operator  $B\varphi$  is defined as  $B(\top, \varphi)$ , which read as “the current agent believes that she satisfies  $\varphi$ ” or “in the most preferred states for the current agent, she satisfies  $\varphi$ .” We can also define the conditional reliability operator  $R_{\mathbf{k}}(\psi, \varphi)$  (read “the most reliable  $\psi$ -agents for agent  $\mathbf{k}$  satisfy  $\varphi$ .”) by

$$R_{\mathbf{k}}(\psi, \varphi) := [1_A](\psi \rightarrow \langle \sqsubseteq_{\mathbf{k}} \rangle (\psi \wedge [\sqsubseteq_{\mathbf{k}}](\psi \rightarrow \varphi))),$$

where we can simplify the clause because of connectedness as noted in (Boutilier, 1994). The unconditional version  $R_{\mathbf{k}}\varphi$  of  $R_{\mathbf{k}}(\psi, \varphi)$  is defined as  $R_{\mathbf{k}}(\top, \varphi)$  which read as “the most reliable agents for agent  $\mathbf{k}$  satisfy  $\varphi$ .” We may also define the “diamond”-version of  $R_{\mathbf{k}}\varphi$  as  $\neg R_{\mathbf{k}}\neg\varphi$  to denote  $\langle R_{\mathbf{k}} \rangle \varphi$ . Then  $\langle R_{\mathbf{k}} \rangle \mathbf{j}$  means that agent  $\mathbf{j}$  is one of the most reliable agents for  $\mathbf{k}$ .

**Example 2.** Let us represent “the current agent likes to put Monet’s picture on wall  $x$ ” by a state-nominal  $\mathbf{a}_x$  in the setting of Example 1. On the *PR*-frame  $F_{exp}$  of Example 1, we define  $V(\mathbf{a}_x) = \{(w_x, i), (w_x, j), (w_x, k)\}$  where  $x = l$  or  $r$ . We use  $\mathbf{i}, \mathbf{j}, \mathbf{k}$  as syntactic names (i.e., agent nominals) for  $i, j$  and  $k$ , where we interpret, e.g.,  $\mathbf{i} = i$  in terms of our valuation function  $V$ . Define  $M_{exp} := (F_{exp}, V)$ . For

example, the preference  $w_l <_i w_r$  can be formalized as a formula  $@_i @_{\mathbf{a}_l} \langle < \rangle \mathbf{a}_r$ , which is valid on  $M_{exp}$ . We can formalize Isabella’s reliability of  $i \approx_i k <_i j$  as  $@_i \langle \sqsubseteq_{\mathbf{i}} \rangle \mathbf{k} \wedge @_{\mathbf{k}} \langle \sqsubseteq_{\mathbf{i}} \rangle \mathbf{i} \wedge @_{\mathbf{k}} \langle \sqsubseteq_{\mathbf{i}} \rangle \mathbf{j}$ , which is valid on  $M_{exp}$ . Moreover,  $@_i \mathbf{Ba}_x$  formalizes “Isabella believes that she likes to put Monet’s picture on wall  $x$ ” and, when  $x = r$ ,  $@_i \mathbf{Ba}_r$  is valid on  $M_{exp}$ . Similarly,  $@_j \mathbf{Ba}_l$  and  $@_{\mathbf{k}} \mathbf{Ba}_r$  are also valid in  $M_{exp}$ . We can see that, from Isabella’s perspective, Ken is one of the most reliable agents who believes that  $\mathbf{a}_r$ . This can be formalized as  $R_{\mathbf{i}}(\mathbf{Ba}_r, \mathbf{k})$ .

The static axiom systems **HPR** and **HPR**<sub>( $m, n$ )</sub> are given as in Table 1, where *uniform substitution* means a substitution that uniformly replaces propositional variables by formulas and nominals from  $\mathbb{N}_i$  by nominals from  $\mathbb{N}_i$  ( $i = 1$  or  $2$ ).

**Theorem 1** (Soundness and completeness).  $\varphi$  is valid in all (possibly infinite) *PR*-models iff  $\varphi$  is derivable in **HPR**. Moreover,  $\varphi$  is valid in all *PR*-models with fixed  $m$  worlds and fixed  $n$  agents iff  $\varphi$  is derivable in **HPR**<sub>( $m, n$ )</sub>. Therefore, **HPR**<sub>( $m, n$ )</sub> is decidable.

We note that the, as far as the authors know, decidability is still unknown for **HPR**, even the fragment of **HPR** without program constructions (cf. (Sano, 2010)). So related computational properties of such fragment has not been yet well-studied (for purely bimodal logic fragment with a slightly different semantics, the reader is referred to (Marx and Mikuláš, 2001)).

### 3 PREFERENCE DYNAMICS

Intuitively, a public announcement of the agents’ individual preferences might induce an agent  $i$  to adjust her own preferences according to what has been announced and the reliability ordering she assigns to the set of agents.<sup>2</sup> For example, an agent might adopt the preferences of the set of agents on whom she relies the most, or might use the strict preferences of her most reliable agents for ‘breaking ties’ among her equally-preferred zones. In (Ghosh and Velázquez-Quesada, 2015a) the authors introduced the *general lexicographic upgrade* operation, which creates a preference ordering following a priority list of orderings. We generalize those operations in the following, where we consider the reliability orderings to be *pre-orders*, rather than being total orders (that is, also anti-symmetric and connected) as they are in the

<sup>2</sup>Note that this work, in line with its predecessor, (Ghosh and Velázquez-Quesada, 2015a), also does not focus on the formal representation of such announcement, but rather on the formal representation of its effects.

Table 1: Axiomatizations HPR and HPR<sub>(m,n)</sub>

<b>Bi-Hybrid Logical Axioms of HPR</b>
All classical tautologies      (Dual <sub>π</sub> ) $\langle \pi \rangle p \leftrightarrow \neg[\pi]\neg p$ (K <sub>π</sub> ) $[\pi](p \rightarrow q) \rightarrow ([\pi]p \rightarrow [\pi]q)$ Let $\mathbf{n} \in N_i$ and $(\mathbf{n}, \mathbf{m}) \in N_i^2$ ( $i = 1, 2$ ) below in this group (K <sub>@</sub> ) $@_{\mathbf{n}}(p \rightarrow q) \rightarrow (@_{\mathbf{n}}p \rightarrow @_{\mathbf{n}}q)$ (SelfDual <sub>@</sub> ) $\neg @_{\mathbf{n}}p \leftrightarrow @_{\mathbf{n}}\neg p$ (Ref) $@_{\mathbf{n}}\mathbf{n}$ (Intro) $\mathbf{n} \wedge p \rightarrow @_{\mathbf{n}}p$ (Agree) $@_{\mathbf{n}}@_{\mathbf{m}}p \rightarrow @_{\mathbf{m}}p$ (Back) $\langle \pi \rangle @_{\mathbf{a}}@_{\mathbf{i}}p \rightarrow @_{\mathbf{a}}@_{\mathbf{i}}p$
<b>Inference Rules of HPR</b>
Modus Ponens, Uniform Substitutions, Necessitation Rules for $[\pi]$ , $@_{\mathbf{i}}$ and $@_{\mathbf{a}}$ (Name) From $\mathbf{n} \rightarrow \varphi$ infer $\varphi$ , where $\mathbf{n} \in N_1 \cup N_2$ is fresh in $\varphi$ (BG <sub>π</sub> ) From $@_{\mathbf{a}}@_{\mathbf{i}}\langle \pi \rangle(\mathbf{b} \wedge \mathbf{j}) \rightarrow @_{\mathbf{b}}@_{\mathbf{j}}\varphi$ infer $@_{\mathbf{a}}@_{\mathbf{i}}[\pi]\varphi$ , where $\mathbf{b}$ and $\mathbf{j}$ are fresh in $@_{\mathbf{a}}@_{\mathbf{i}}[\pi]\varphi$
<b>Interaction Axioms of HPR</b>
(Com@) $@_{\mathbf{i}}@_{\mathbf{a}}p \leftrightarrow @_{\mathbf{a}}@_{\mathbf{i}}p$ (Red@ <sub>1</sub> ) $\mathbf{a} \leftrightarrow @_{\mathbf{i}}\mathbf{a}$ (Red@ <sub>2</sub> ) $\mathbf{i} \leftrightarrow @_{\mathbf{a}}\mathbf{i}$ (Dcom <sub>W</sub> @ <sub>2</sub> ) $@_{\mathbf{i}}\langle \alpha \rangle p \leftrightarrow @_{\mathbf{i}}\langle \alpha \rangle @_{\mathbf{i}}p$ ( $\alpha \in \{1_W, \leq, \geq\}$ ) (Com <sub>A</sub> @ <sub>1</sub> ) $@_{\mathbf{a}}\langle \beta \rangle p \leftrightarrow \langle \beta \rangle @_{\mathbf{a}}p$ ( $\beta \in \{1_A, \sqsubseteq_{\mathbf{k}}, \sqsupset_{\mathbf{k}}\}$ )
<b>Axioms for Atomic Programs of HPR</b>
(U <sub>W</sub> ) $@_{\mathbf{a}}\langle 1_W \rangle \mathbf{b}$ (Cnv <sub>≤</sub> ) $@_{\mathbf{a}}\langle \leq \rangle \mathbf{b} \leftrightarrow @_{\mathbf{b}}\langle \geq \rangle \mathbf{a}$ (U <sub>A</sub> ) $@_{\mathbf{i}}\langle 1_A \rangle \mathbf{j}$ (Cnv <sub>⊆</sub> ) $@_{\mathbf{i}}\langle \sqsubseteq_{\mathbf{k}} \rangle \mathbf{j} \leftrightarrow @_{\mathbf{j}}\langle \sqsupset_{\mathbf{k}} \rangle \mathbf{i}$ (Eq <sub>⊆</sub> ) $@_{\mathbf{i}}\mathbf{j} \rightarrow ([\sqsubseteq_{\mathbf{i}}]p \rightarrow [\sqsubseteq_{\mathbf{j}}]p)$
<b>Axioms for Compounded Programs of HPR</b>
(∪) $\langle \pi \cup \rho \rangle p \leftrightarrow \langle \pi \rangle p \vee \langle \rho \rangle p$ (?) $\langle ?(\varphi, \psi) \rangle p \leftrightarrow \varphi \wedge \langle 1_A \rangle \langle 1_W \rangle (\psi \wedge p)$ (∩) $@_{\mathbf{a}}@_{\mathbf{i}}\langle \pi \cap \rho \rangle (\mathbf{b} \wedge \mathbf{j}) \leftrightarrow @_{\mathbf{a}}@_{\mathbf{i}}(\langle \pi \rangle (\mathbf{b} \wedge \mathbf{j}) \wedge \langle \rho \rangle (\mathbf{b} \wedge \mathbf{j}))$ (− <sub>W</sub> ) $@_{\mathbf{a}}\langle -\alpha \rangle \mathbf{b} \leftrightarrow @_{\mathbf{a}}\neg \langle \alpha \rangle \mathbf{b}$ ( $\alpha \in \{1_W, \leq, \geq\}$ ) (− <sub>A</sub> ) $@_{\mathbf{i}}\langle -\beta \rangle \mathbf{j} \leftrightarrow @_{\mathbf{i}}\neg \langle \beta \rangle \mathbf{j}$ ( $\beta \in \{1_A, \sqsubseteq_{\mathbf{k}}, \sqsupset_{\mathbf{k}}\}$ ) (∏ <sub>i</sub> ) $@_{\mathbf{a}}@_{\mathbf{k}}\langle (\pi, \mathbf{j}) \prod_{\mathbf{i}} (\pi', \mathbf{j}') \rangle (\mathbf{b} \wedge \mathbf{k}') \leftrightarrow @_{\mathbf{i}}(\mathbf{k} \wedge \mathbf{k}') \wedge @_{\mathbf{a}}@_{\mathbf{j}}\langle \pi \rangle (\mathbf{b} \wedge \mathbf{j}) \wedge @_{\mathbf{a}}@_{\mathbf{j}'}\langle \pi' \rangle (\mathbf{b} \wedge \mathbf{j}')$
<b>Axioms for PR-frames of HPR</b>
(4 <sub>≤</sub> ) $@_{\mathbf{a}}\langle \leq \rangle \mathbf{b} \wedge @_{\mathbf{b}}\langle \leq \rangle \mathbf{c} \rightarrow @_{\mathbf{a}}\langle \leq \rangle \mathbf{c}$ (4 <sub>⊆</sub> ) $@_{\mathbf{j}}\langle \sqsubseteq_{\mathbf{i}} \rangle \mathbf{k} \wedge @_{\mathbf{k}}\langle \sqsubseteq_{\mathbf{i}} \rangle \mathbf{l} \rightarrow @_{\mathbf{j}}\langle \sqsubseteq_{\mathbf{i}} \rangle \mathbf{l}$ (Ref <sub>≤</sub> ) $@_{\mathbf{a}}\langle \leq \rangle \mathbf{a}$ (Cmp <sub>⊆</sub> ) $@_{\mathbf{j}}\langle \sqsubseteq_{\mathbf{i}} \rangle \mathbf{k} \vee @_{\mathbf{k}}\langle \sqsubseteq_{\mathbf{i}} \rangle \mathbf{j}$
<b>Additional Axioms and Rules for HPR<sub>(m,n)</sub></b>
( W  ≤ m) $\bigvee_{0 \leq k \neq l \leq m} @_{\mathbf{a}_k} \mathbf{a}_l$ ( A  ≤ n) $\bigvee_{0 \leq k \neq l \leq n} @_{\mathbf{i}_k} \mathbf{i}_l$ ( W  ≥ m) From $(\bigwedge_{1 \leq k \neq l \leq m} \neg @_{\mathbf{a}_k} \mathbf{a}_l) \rightarrow \psi$ infer $\psi$ , where $\mathbf{a}_k$ s are fresh in $\psi$ ( A  ≥ n) From $(\bigwedge_{1 \leq k \neq l \leq n} \neg @_{\mathbf{i}_k} \mathbf{i}_l) \rightarrow \psi$ infer $\psi$ , where $\mathbf{i}_k$ s are fresh in $\psi$ .

earlier work, which was quite an artificial assumption on agents' reliabilities. Agent  $i$ 's preference ordering *after* an announcement,  $\leq'_i$ , can be defined in terms of the just announced preferences (the agents' preferences *before* the announcement,  $\leq_1, \dots, \leq_n$ ) and how much  $i$  relied on each agent ( $i$ 's reliability *before* the announcement,  $\leq_i$ ):  $\leq'_i := f(\leq_1, \dots, \leq_n, \leq_i)$  for some function  $f$ . Here are some such functions inspired by (van Benthem, 2007; Ghosh and Velázquez-Quesada, 2015a).

**Definition 5.** Given a set  $X \subseteq A$  of agents,  $u <_X v$  if  $u <_k v$  holds for all agents  $k \in X$ . Moreover,  $u \succ_X v$  is used to mean  $u <_X v$  or  $v <_X u$  and  $\text{dom}(\succ_X) := \{u \in A \mid u \succ_X v \text{ for some } v \in A\}$ .

Note that  $\text{dom}(\succ_X)$  allows us to specify the connected components by the relation  $\succ_X$ . Recall that  $mr(i)$  denotes the set of all maximally reliable agents for  $i$ .

**Definition 6** (Conservative Upgrade). Agent  $i$  takes the strict preference ordering of her most reliable agents, and leaves the rest undecided (equipreferable). More precisely, the upgraded ordering  $\leq'_i$  is defined by:  $u \leq'_i v$  iff ( $u <_{mr(i)} v$  or  $u = v$ ) or ( $u, v \notin \text{dom}(\succ_X)$ ).

**Definition 7** (Radical Upgrade). Agent  $i$  takes the strict preference ordering of her most reliable agents, and in the remaining disjoint zones she uses her old ordering. More precisely, the upgraded ordering  $\leq'_i$  is defined by:  $u \leq'_i v$  iff ( $u <_{mr(i)} v$  or  $u = v$ ) or ( $u, v \notin \text{dom}(\succ_X)$  and  $u \leq_i v$ ).

Note that both the conservative and radical upgrades preserve preorders (and thus upgraded models belong to our class of semantic models).

### 3.1 Expressing the preference dynamics

To formalize preference dynamics from the previous section, we add the following dynamic operators to the static syntax  $\mathcal{HL}$ . First of all, we regard all the agents involved in our two preference upgrade above as agent nominals (syntactic names of agents) and so let us denote agent  $i$ 's syntactic name as  $\mathbf{i}$  of boldface and the set of all syntactic names in  $mr(i)$  as  $\mathbf{mr}(\mathbf{i})$ .  $\mathcal{HL}_{\{\text{pu}\}}$  is defined to be an expansion of  $\mathcal{HL}$  with all operators  $\langle \text{pu}_{\mathcal{R}}^{\mathbf{i}} \rangle$ , where  $\mathbf{i}$  be an agent-nominal and  $\mathcal{R}$  is a list of sets of agent-nominals defined as  $\mathcal{R} = \mathbf{mr}(\mathbf{i})$  (conservative upgrade) or  $\mathcal{R} = (\mathbf{mr}(\mathbf{i}); \{\mathbf{i}\})$  (radical upgrade).

**Definition 8** (Operators). A formula  $\text{Req}(\mathcal{R})$ , representing requirements for the list  $\mathcal{R}$  is defined as the conjunction  $\bigwedge_{\mathbf{j} \neq \mathbf{k} \in \mathbf{mr}(\mathbf{i})} \neg @_{\mathbf{j}} \mathbf{k}$  (mutual disjointness of agents involved in  $mr(i)$ ) and  $\bigwedge_{\mathbf{j} \in \mathbf{mr}(\mathbf{i})} \langle R_{\mathbf{i}} \rangle \mathbf{j}$  ( $mr(i)$  is the set of maximally reliable agents for  $\mathbf{i}$ ). Given a

$PR$ -model  $M = (W, \{\leq_i, \leq_i\}_{i \in A}, V)$ , define:

$$M, (w, j) \Vdash \langle \text{pu}_{\mathcal{R}}^{\mathbf{i}} \rangle \varphi \quad \text{iff} \quad M, (w, j) \Vdash \text{Req}(\mathcal{R})$$

$$\text{and } \text{pu}_{\mathcal{R}}^{\mathbf{i}}(M), (w, j) \Vdash \varphi,$$

where  $\text{pu}_{\mathcal{R}}^{\mathbf{i}}(M)$  is the same model as  $M$  except  $\leq_i$  is replaced by  $\leq_{\mathcal{R}}$  where  $\mathcal{R} = \mathbf{mr}(\mathbf{i})$  or  $\mathcal{R} = (\mathbf{mr}(\mathbf{i}); \{\mathbf{i}\})$  and corresponding  $\leq_{\mathcal{R}}$ 's are given by Definitions 6 and 7, respectively.

For an axiom system for the modality  $\langle \text{pu}_{\mathcal{R}}^{\mathbf{i}} \rangle$ , we will provide recursion axioms: valid formulas and validity-preserving rules indicating how to translate a formula with the new modality into a provably equivalent one without them. In this case, the modalities can take the form of any relational expression. So we provide a ‘matching’ relational expression in the original model  $M$  by defining relational transformers similar to those in (Ghosh and Velázquez-Quesada, 2015a; Ghosh and Velázquez-Quesada, 2015b), in spirit of the program transformers of (van Benthem et al., 2006).

Before going into the notion of relational transformer, we have two observations. Firstly, when  $\pi := ?\mathbf{j} \cap \leq$ , we note that  $(w, i)R_{\pi}(v, k)$  is equivalent to the conjunction of  $i = k = \mathbf{j}$  and  $w \leq_{\mathbf{j}} v$ . Similarly, when we put  $\pi' := ?\neg \mathbf{j} \cap \leq$ , we remark that  $(w, i)R_{\pi'}(v, k)$  is equivalent to the conjunction of  $w \leq_i v$  and  $i = k$  and  $i \neq \mathbf{j}$ . Secondly, to reflect the relation  $<_X$  in Definition 5, we need our program construction  $(\pi, \mathbf{j}) \sqcap_{\mathbf{i}} (\rho, \mathbf{k})$  to taking the intersection of (strict) preference relations of the possibly different agents than  $\mathbf{i}$ . These observations allow us to capture the idea behind Definitions 6 and 7 syntactically in the following definition.

**Definition 9** (Relational transformer). Let us introduce the following abbreviations for relational expressions: We define  $<_{\mathbf{mr}(\mathbf{i})} := \sqcap_{\mathbf{i}} \{(\leq \cap \neg \geq, \mathbf{j}) \mid \mathbf{j} \in \mathbf{mr}(\mathbf{i})\}$ . Then  $>_{\mathbf{mr}(\mathbf{i})}$  is similarly defined and  $\asymp_{\mathbf{mr}(\mathbf{i})}$  is defined to be  $<_{\mathbf{mr}(\mathbf{i})} \cup >_{\mathbf{mr}(\mathbf{i})}$ . Moreover, a formula  $\mathbf{d}(\asymp_{\mathbf{mr}(\mathbf{i})})$  is defined as  $\langle \asymp_{\mathbf{mr}(\mathbf{i})} \rangle \top$ .

A relational transformer  $Tu_{\mathcal{R}}^{\mathbf{i}}$  is a function from relational expressions to relational expressions defined as follows. When  $\mathcal{R} = \mathbf{mr}(\mathbf{i})$  (conservative upgrade),

$$Tu_{\mathcal{R}}^{\mathbf{i}}(\alpha) := \alpha \quad (\alpha \in \{1_A, 1_W, \sqsubseteq_{\mathbf{k}}, \sqsupseteq_{\mathbf{k}} \mid \mathbf{k} \in N_2\}),$$

$$Tu_{\mathcal{R}}^{\mathbf{i}}(\leq) := \left( ?\mathbf{i} \cap (<_{\mathbf{mr}(\mathbf{i})} \cup 1_A \cup ?\neg \mathbf{d}(\asymp_{\mathbf{mr}(\mathbf{i})})) \right) \cup (? \neg \mathbf{i} \cap \leq)$$

$$Tu_{\mathcal{R}}^{\mathbf{i}}(\geq) := \left( ?\mathbf{i} \cap (>_{\mathbf{mr}(\mathbf{i})} \cup 1_A \cup ?\neg \mathbf{d}(\asymp_{\mathbf{mr}(\mathbf{i})})) \right) \cup (? \neg \mathbf{i} \cap \geq)$$

$$Tu_{\mathcal{R}}^{\mathbf{i}}(\pi \cup \rho) := Tu_{\mathcal{R}}^{\mathbf{i}}(\pi) \cup Tu_{\mathcal{R}}^{\mathbf{i}}(\rho),$$

$$Tu_{\mathcal{R}}^{\mathbf{i}}(\pi \cap \rho) := Tu_{\mathcal{R}}^{\mathbf{i}}(\pi) \cap Tu_{\mathcal{R}}^{\mathbf{i}}(\rho),$$

$$Tu_{\mathcal{R}}^{\mathbf{i}}(?(\varphi, \psi)) := ?(\langle \text{pu}_{\mathcal{R}}^{\mathbf{i}} \rangle \varphi, \langle \text{pu}_{\mathcal{R}}^{\mathbf{i}} \rangle \psi).$$

$$Tu_{\mathcal{R}}^{\mathbf{i}}((\pi, \mathbf{k}) \sqcap_{\mathbf{j}} (\rho, \mathbf{l})) := ((Tu_{\mathcal{R}}^{\mathbf{i}}(\pi), \mathbf{k}) \sqcap_{\mathbf{j}} (Tu_{\mathcal{R}}^{\mathbf{i}}(\rho), \mathbf{l}))$$

$$Tu_{\mathcal{R}}^{\mathbf{i}}(-\beta) := -Tu_{\mathcal{R}}^{\mathbf{i}}(\beta),$$

where  $\beta \in \{1_W, \leq, \geq\} \cup \{1_A, \sqsubseteq_{\mathbf{k}}, \sqsupseteq_{\mathbf{k}} \mid \mathbf{k} \in N_2\}$ . When  $\mathcal{R} = (\mathbf{mr}(\mathbf{i}); \{\mathbf{i}\})$ , we replace the occurrence of “ $? \neg \mathbf{d}(\asymp_{\mathbf{mr}(\mathbf{i})})$ ” in  $Tu_{\mathcal{R}}^{\mathbf{i}}(\leq)$  or  $Tu_{\mathcal{R}}^{\mathbf{i}}(\geq)$  with

$$“? \neg \mathbf{d}(\asymp_{\mathbf{mr}(\mathbf{i})}) \top \cap \leq” \quad \text{or} \quad “? \neg \mathbf{d}(\asymp_{\mathbf{mr}(\mathbf{i})}) \cap \leq,”$$

respectively.

**Theorem 2.** The axioms and rules below together with those of **HPR** (or, those of **HPR**<sub>(m,n)</sub>) provide sound and complete axiom systems for  $\mathcal{HL}_{\{\text{pu}\}}$  with respect to possibly infinite  $PR$  models (or,  $PR$  models with  $m$  worlds and  $n$  agents, respectively).

$$\langle \text{pu}_{\mathcal{R}}^{\mathbf{i}} \rangle p \leftrightarrow \text{Req}(\mathcal{R}) \wedge p,$$

$$\langle \text{pu}_{\mathcal{R}}^{\mathbf{i}} \rangle (\varphi \vee \psi) \leftrightarrow \langle \text{pu}_{\mathcal{R}}^{\mathbf{i}} \rangle \varphi \vee \langle \text{pu}_{\mathcal{R}}^{\mathbf{i}} \rangle \psi,$$

$$\langle \text{pu}_{\mathcal{R}}^{\mathbf{i}} \rangle \neg \varphi \leftrightarrow \text{Req}(\mathcal{R}) \wedge \neg \langle \text{pu}_{\mathcal{R}}^{\mathbf{i}} \rangle \varphi$$

$$\langle \text{pu}_{\mathcal{R}}^{\mathbf{i}} \rangle \mathbf{j} \leftrightarrow \text{Req}(\mathcal{R}) \wedge \mathbf{j}, \quad \langle \text{pu}_{\mathcal{R}}^{\mathbf{i}} \rangle \mathbf{a} \leftrightarrow \text{Req}(\mathcal{R}) \wedge \mathbf{a},$$

$$\langle \text{pu}_{\mathcal{R}}^{\mathbf{i}} \rangle @_{\mathbf{j}} \varphi \leftrightarrow \text{Req}(\mathcal{R}) \wedge @_{\mathbf{j}} \langle \text{pu}_{\mathcal{R}}^{\mathbf{i}} \rangle \varphi,$$

$$\langle \text{pu}_{\mathcal{R}}^{\mathbf{i}} \rangle @_{\mathbf{a}} \varphi \leftrightarrow \text{Req}(\mathcal{R}) \wedge @_{\mathbf{a}} \langle \text{pu}_{\mathcal{R}}^{\mathbf{i}} \rangle \varphi$$

$$\langle \text{pu}_{\mathcal{R}}^{\mathbf{i}} \rangle \langle \pi \rangle \varphi \leftrightarrow \text{Req}(\mathcal{R}) \wedge \langle Tu_{\mathcal{R}}^{\mathbf{i}}(\pi) \rangle \langle \text{pu}_{\mathcal{R}}^{\mathbf{i}} \rangle \varphi,$$

From  $\varphi \rightarrow \psi$ , we may infer  $\langle \text{pu}_{\mathcal{R}}^{\mathbf{i}} \rangle \varphi \rightarrow \langle \text{pu}_{\mathcal{R}}^{\mathbf{i}} \rangle \psi$ .

*Proof.* Soundness of the new axioms are straightforward. Completeness follows from the completeness of the static system **HPR** (cf. Chapter 7 of (van Ditmarsch et al., 2008), for an extensive explanation of this technique).  $\square$

**Example 3.** In our running example of Section 1, each agent is regarded to employ conservative upgrades to change his or her preference. Let us write the corresponding upgrade operators of  $\{i, j, k\}$  by  $\langle \text{pu}_{\mathcal{R}_i}^{\mathbf{i}} \rangle$  and  $\langle \text{pu}_{\mathcal{R}_j}^{\mathbf{j}} \rangle$ ,  $\langle \text{pu}_{\mathcal{R}_k}^{\mathbf{k}} \rangle$ , respectively. Then, three flatmates did not reach an agreement after conservative upgrades of all agents, i.e.,

$$@_{\mathbf{i}} \mathbf{Ba}_r \wedge @_{\mathbf{j}} \mathbf{Ba}_l \wedge @_{\mathbf{k}} \mathbf{Ba}_r \wedge$$

$$\langle \text{pu}_{\mathcal{R}_i}^{\mathbf{i}} \rangle \langle \text{pu}_{\mathcal{R}_j}^{\mathbf{j}} \rangle \langle \text{pu}_{\mathcal{R}_k}^{\mathbf{k}} \rangle (@_{\mathbf{i}} \mathbf{Ba}_l \wedge @_{\mathbf{j}} \mathbf{Ba}_r \wedge @_{\mathbf{k}} \mathbf{Ba}_r).$$

is valid in  $M_{exp}$ , because upgraded preferences are given by  $w_r <'_i w_l$ ,  $w_l <'_j w_r$  and  $w_l <'_k w_r$ .

## 4 RELIABILITY DYNAMICS

A public announcement of the agents’ individual preferences may change the agents’ reliability attributions as well: for example, an agent may consider more reliable those agents whose preferences coincide (or, for some reason, differ) from her own. In such cases, agent  $i$ ’s new reliability ordering,  $\leq'_i$ ,

can be given in terms of the agents' current preferences,  $\leq_1, \dots, \leq_n$ , and  $i$ 's current reliability ordering,  $\leq_i$ . Thus,  $\leq'_i := g(\leq_1, \dots, \leq_n, \leq_i)$  for some function  $g$ . We now provide formal definitions of some such possibilities.

#### 4.1 Reliability change operations

The notion of ‘‘matching preference orders’’ will form the basis for the reliability dynamics. The idea is that two preference orderings match each other to a certain extent if the orderings are identical on some subset of the state space. A full match indicates that the orderings coincide on the whole domain; a partial match indicates that they coincide up to some proper subset of the domain.

**Definition 10** (Matching preferences). Let  $F$  be a  $PR$  frame given by  $F = (W, \{\leq_i, \leq'_i\}_{i \in A})$  and let  $i \in A$  be an agent. If  $\leq_i$  is identical with  $\leq_j$  on  $W' \subseteq W$ , then  $W'$  is said to be a set of match for  $i$  and  $j$  (notation:  $\leq_i \sim_{W'} \leq_j$ ).

- Preference orders  $\leq_i$  and  $\leq_j$  are said to *fully match* each other iff  $\leq_i \sim_W \leq_j$ .<sup>3</sup>  $FullMat(i)$  denotes the set of agents in  $A \setminus \{i\}$  having full match with  $i$ .
- Preference orders  $\leq_i$  and  $\leq_j$  have *zero match* with each other iff there is no  $W' \subseteq W$  with  $|W'| \geq 2$  such that  $\leq_i \sim_{W'} \leq_j$ .<sup>4</sup>  $ZeroMat(i)$  denotes the set of agents in  $A \setminus \{i\}$  having zero match with  $i$ .

With these definitions, we can define some operations for reliability change.

**Definition 11** (Full, Zero matching upgrade). Agent  $i$  puts those agents that have full/zero match with her own preference ordering above those that do not, keeping her old reliability ordering within each of the two zones. More precisely, if  $\leq_i$  is agent  $i$ 's current reliability ordering, then her new reliability ordering  $\leq'_i$  is defined by:

$$j \leq'_i k \text{ iff } (j, k \in V \text{ and } j \leq_i k) \text{ or } (k \in V \text{ and } j \notin V) \\ \text{or } (j, k \notin V \text{ and } j \leq_i k).$$

Here  $V = FullMat(i) \cup \{i\}, ZeroMat(i)$ , respectively.

Once again, we can consider more generalized definitions for upgrade policies as well, but we just stick to simple definitions to give the main idea. Note that both the full matching and zero matching upgrades preserve total preorders (and thus upgraded models belong to our class of semantic models).

<sup>3</sup>Note how, by the finiteness of  $W$  (the reflexivity of the preference relations), there is always a maximal  $X \subseteq W$  such that  $\leq_i \sim_X \leq_j$  for every agent  $i, j$ .

<sup>4</sup>For the same reason, there is always a minimal  $X \subseteq W$  such that  $\leq_i \sim_X \leq_j$  for every agent  $i, j$ .

#### 4.2 Expressing the reliability dynamics

To describe reliability dynamics from the previous section, the following dynamic operators are added to the static syntax of  $\mathcal{HL}$ .  $\mathcal{HL}_{\{rc\}}$  is defined to be an expansion of  $\mathcal{HL}$  with all operators  $\langle rc_{\mathcal{E}}^i \rangle$ , where  $\mathbf{i}$  is an agent-nominal and  $\mathcal{E}$  a pair of  $\mathcal{HL}$ -formulas of the form  $@_{\mathbf{a}}\chi$  (recall:  $\mathbf{a}$  is a world-nominal). An underlying semantic intuition for  $\langle rc_{\mathcal{E}}^i \rangle$  is: Given a  $PR$ -model  $M$ , the pair  $\mathcal{E} = (@_{\mathbf{a}_1}\chi_1, @_{\mathbf{a}_2}\chi_2)$  can be regarded as a partition (i.e., an equivalence relation on agents) in the sense that  $\{\{i \in A \mid M, (i, \mathbf{a}_k) \models \chi_k\}\}_{1 \leq k \leq 2}$  forms a partition of  $A$ , and the reliability ordering  $\leq_i$  of the original  $PR$  model  $M$  is rewritten into the updated reliability ordering  $\leq'_i$  as in Definition 11 of the Section 4.1.

**Definition 12** (Operators). Given any pair  $\mathcal{E} = (\varphi_1, \varphi_2)$  of formulas of the form  $@_{\mathbf{a}}\chi$ , a formula  $Eq(\mathcal{E})$  is defined as the conjunction of  $[1_A](\varphi_1 \vee \varphi_2)$  (exhaustiveness for agents) and  $\neg(1_A)(\varphi_1 \wedge \varphi_2)$  (pairwise disjointness for agents). Given a pair  $\mathcal{E} = (@_{\mathbf{a}_1}\chi_1, @_{\mathbf{a}_2}\chi_2)$  and a  $PR$ -model  $M = (W, \{\leq_i, \leq'_i\}_{i \in A}, V)$ , define:

$$M, (w, j) \models \langle rc_{\mathcal{E}}^i \rangle \varphi \text{ iff } M, (w, j) \models Eq(\mathcal{E}) \\ \text{and } rc_{\mathcal{E}}^i(M), (w, j) \models \varphi,$$

where  $rc_{\mathcal{E}}^i(M)$  is the same model as  $M$  except  $\leq_i$  is replaced by  $\leq'_i$  of Definition 11.

**Definition 13** (Relational transformer). Let  $\mathcal{E} = (\varphi_1, \varphi_2)$  be a pair. A relational transformer  $Tr_{\mathcal{E}}^i$  is a function from relational expressions to relational expressions defined as follows.

$$Tr_{\mathcal{E}}^i(\alpha) := \alpha \quad (\alpha \in \{1_A, 1_W, \leq, \geq\}), \\ Tr_{\mathcal{E}}^i(\sqsubseteq_i) := (\sqsubseteq_i \cap (? \varphi_1 \cup ? \varphi_2)) \cup (1_A \cap (? \varphi_1, \varphi_2)), \\ Tr_{\mathcal{E}}^i(\sqsupseteq_i) := (\sqsupseteq_i \cap (? \varphi_1 \cup ? \varphi_2)) \cup (1_A \cap (? \varphi_1, \varphi_2)), \\ Tr_{\mathcal{E}}^i(\sqsubseteq_{\mathbf{k}}) := (? @_{\mathbf{i}} \mathbf{k} \cap Tr_{\mathcal{E}}^i(\sqsubseteq_i)) \cup (? \neg @_{\mathbf{i}} \mathbf{k} \cap \sqsubseteq_{\mathbf{k}}) \quad (\mathbf{k} \neq \mathbf{i}), \\ Tr_{\mathcal{E}}^i(\sqsupseteq_{\mathbf{k}}) := (? @_{\mathbf{i}} \mathbf{k} \cap Tr_{\mathcal{E}}^i(\sqsupseteq_i)) \cup (? \neg @_{\mathbf{i}} \mathbf{k} \cap \sqsupseteq_{\mathbf{k}}) \quad (\mathbf{k} \neq \mathbf{i}), \\ Tr_{\mathcal{E}}^i(\pi \cup \rho) := Tr_{\mathcal{E}}^i(\pi) \cup Tr_{\mathcal{E}}^i(\rho), \\ Tr_{\mathcal{E}}^i(\pi \cap \rho) := Tr_{\mathcal{E}}^i(\pi) \cap Tr_{\mathcal{E}}^i(\rho), \\ Tr_{\mathcal{E}}^i(?(\varphi, \psi)) := ?(\langle rc_{\mathcal{E}}^i \rangle \varphi, \langle rc_{\mathcal{E}}^i \rangle \psi). \\ Tr_{\mathcal{E}}^i((\pi, \mathbf{k}) \sqcap_j (\rho, \mathbf{k})) := (Tr_{\mathcal{E}}^i(\pi), \mathbf{k}) \sqcap_j (Tr_{\mathcal{E}}^i(\rho), \mathbf{k}), \\ Tr_{\mathcal{E}}^i(-\alpha) := -Tr_{\mathcal{E}}^i(\alpha),$$

where  $\alpha \in \{1_W, \leq, \geq\} \cup \{1_A, \sqsubseteq_{\mathbf{k}}, \sqsupseteq_{\mathbf{k}} \mid \mathbf{k} \in N_2\}$ .

When  $\mathbf{k} \neq \mathbf{i}$ , i.e.,  $\mathbf{k}$  and  $\mathbf{i}$  are syntactically distinct agent nominals, the reader may wonder why we should have generalized test operators ‘‘ $? @_{\mathbf{i}} \mathbf{k}$ ’’ and ‘‘ $? \neg @_{\mathbf{i}} \mathbf{k}$ ’’ in the definitions  $Tr_{\mathcal{E}}^i(\sqsubseteq_{\mathbf{k}})$  and  $Tr_{\mathcal{E}}^i(\sqsupseteq_{\mathbf{k}})$ . This is because the

same agent might have two distinct (syntactic) names. Based on a similar strategy for Theorem 2, we can now prove the following theorem.

**Theorem 3.** The axioms and rules below together with those of **HPR** (or, those of **HPR**<sub>(m,n)</sub>) provide sound and complete axiom systems for  $\mathcal{HL}_{\{rc\}}$  with respect to possibly infinite *PR* models (or, *PR* models with  $m$  worlds and  $n$  agents, respectively).

$$\begin{aligned} \langle rc_{\mathcal{E}}^i \rangle p &\leftrightarrow \text{Eq}(\mathcal{E}) \wedge p, \quad \langle rc_{\mathcal{E}}^i \rangle (\varphi \vee \psi) \leftrightarrow \langle rc_{\mathcal{E}}^i \rangle \varphi \vee \langle rc_{\mathcal{E}}^i \rangle \psi, \\ \langle rc_{\mathcal{E}}^i \rangle \neg \varphi &\leftrightarrow \text{Eq}(\mathcal{E}) \wedge \neg \langle rc_{\mathcal{E}}^i \rangle \varphi \\ \langle rc_{\mathcal{E}}^i \rangle \mathbf{j} &\leftrightarrow \text{Eq}(\mathcal{E}) \wedge \mathbf{j}, \quad \langle rc_{\mathcal{E}}^i \rangle \mathbf{a} \leftrightarrow \text{Eq}(\mathcal{E}) \wedge \mathbf{a}, \\ \langle rc_{\mathcal{E}}^i \rangle @_j \varphi &\leftrightarrow \text{Eq}(\mathcal{E}) \wedge @_j \langle rc_{\mathcal{E}}^i \rangle \varphi, \\ \langle rc_{\mathcal{E}}^i \rangle @_a \varphi &\leftrightarrow \text{Eq}(\mathcal{E}) \wedge @_a \langle rc_{\mathcal{E}}^i \rangle \varphi \\ \langle rc_{\mathcal{E}}^i \rangle \langle \pi \rangle \varphi &\leftrightarrow \text{Eq}(\mathcal{E}) \wedge \langle Tr_{\mathcal{E}}^i(\pi) \rangle \langle rc_{\mathcal{E}}^i \rangle \varphi, \\ \text{From } \varphi \rightarrow \psi, &\text{ we may infer } \langle rc_{\mathcal{E}}^i \rangle \varphi \rightarrow \langle rc_{\mathcal{E}}^i \rangle \psi. \end{aligned}$$

**Example 4.** After Isabella and John know others' preferences, we regard, in our running example, that Isabella uses full-match reliability change  $\langle rc_{\mathcal{E}_i}^i \rangle$  and John employs zero-match reliability change  $\langle rc_{\mathcal{E}_j}^j \rangle$ . Unlike Example 3, let us first consider reliability changes of Isabella and John and then take the conservative upgrades of all agents. This process and the resulting agreements among agents are formalized as

$$\begin{aligned} &\langle rc_{\mathcal{E}_i}^i \rangle \langle rc_{\mathcal{E}_j}^j \rangle \langle pu_{\mathcal{R}_i}^i \rangle \langle pu_{\mathcal{R}_j}^j \rangle \langle pu_{\mathcal{R}_k}^k \rangle \\ &(@_i Ba_r \wedge @_j Ba_r \wedge @_k Ba_r), \end{aligned}$$

which is valid in  $M_{exp}$ , because Isabella's reliability is changed into  $j <'_i i \approx'_i k$  and John's reliability does not change.

We note here that while the main focus of the work is to model joint deliberation in form of simultaneous preference and reliability upgrades, the model operations and modalities of Sections 3.1 and 4.2 deal with single agent upgrades. This presentation style has been chosen in order to simplify notation and readability, but the provided definitions can be easily extended in order to match our goals. In particular, the model operations of Definitions 8 and 12 can be extended to simultaneous upgrades by asking for a list  $\mathcal{R}$  of lexicographic lists (with  $\mathcal{R}_i$  the list for agent  $i$ ), and asking for a list  $\mathcal{E}$  of partition lists (with  $\mathcal{E}_i$  the list for agent  $i$ ), respectively. Then the corresponding modalities,  $\langle pu_{\mathcal{R}}^i \rangle$  and  $\langle rc_{\mathcal{E}}^i \rangle$  can still be axiomatised by the presented system with some simple modifications.

## 5 CONCLUSION

This work continues the line of study in (Ghosh and Velázquez-Quesada, 2015a; Ghosh and Velázquez-Quesada, 2015b) and provides a further interplay between the preferences that the agents have about the world around and the reliability attributions they have with respect to one another. We deal with both preference change based on reliability, and reliability change based on preferences, and propose two-dimensional dynamic hybrid logics to express such changes. The main technical results that we have are sound and complete axiomatizations which lead to decidability (provided the numbers of agents and of states are *fixed* finite numbers) as well. In process, we also discuss about agent beliefs in such situations, e.g. relating reliability attributions with the notions of belief (cf. the running example in the text). The novel contribution of the work is the study of change in reliability attribution of agents based on their preferences.

To conclude, let us provide some pointers towards future work: (1) What other reasonable preference and reliability upgrade policies can there be and how to model them? (2) How to investigate the role of knowledge in such changes, especially if manipulation comes into play? (3) What would be the characterizing conditions for reaching consensus in such deliberative processes? We endeavor to provide answers to such questions in future.<sup>5</sup>

## REFERENCES

- Arrow, K. J., Sen, A. K., and Suzumura, K., editors (2002). *Handbook of Social Choice and Welfare*. Elsevier. Two volumes.
- Boutilier, C. (1994). Conditional logics of normality: A modal approach. *Artificial Intelligence*, 68(1):87–154.
- Burgess, J. P. (1984). Basic tense logic. In Gabbay, D. and Guenther, F., editors, *Handbook of Philosophical Logic*, volume II, chapter 2, pages 89–133. Reidel.
- Demolombe, R. (2001). To trust information sources: A proposal for a modal logic framework. In Castelfranchi, C. and Tan, Y.-H., editors, *Trust*

<sup>5</sup>The authors would like to thank the anonymous reviewers for their helpful and constructive comments that greatly contributed to improving the final version of the paper. The work of the second author was partially supported by JSPS KAKENHI Grant-in-Aid for Young Scientists (B) Grant Number 15K21025 and JSPS Core-to-Core Program (A. Advanced Research Networks).



- and *Deception in Virtual Societies*. Kluwer Academic, Dordrecht.
- Demolombe, R. (2004). Reasoning about trust: A formal logical framework. In Jensen, C. D., Poslad, S., and Dimitrakos, T., editors, *iTrust*, volume 2995 of *Lecture Notes in Computer Science*, pages 291–303. Springer.
- Endriss, U. (2011). Logic and social choice theory. In Gupta, A. and van Benthem, J., editors, *Logic and Philosophy Today*, volume 2, pages 333–377. College Publications.
- Falcone, R., Barber, K. S., Sabater-Mir, J., and Singh, M. P., editors (2008). *Trust in Agent Societies, 11th International Workshop, TRUST 2008, Estoril, Portugal, May 12-13, 2008. Revised Selected and Invited Papers*, volume 5396 of *Lecture Notes in Computer Science*. Springer.
- Falcone, R. and Castelfranchi, C. (2001). Social trust: A cognitive approach. In Castelfranchi, C. and Tan, Y.-H., editors, *Trust and Deception in Virtual Societies*, pages 55–90. Kluwer Academic, Dordrecht.
- Gargov, G., Passy, S., and Tinchev, T. (1987). Modal environment for boolean speculations, preliminary report. In Skordev, D., editor, *Mathematical Logic and Its Applications*, pages 253–263. Plenum Press.
- Ghosh, S. and Velázquez-Quesada, F. R. (2015a). Agreeing to agree: Reaching unanimity via preference dynamics based on reliable agents. In Bordini, R., Elkind, E., Weiss, G., and Yolum, P., editors, *AAMAS 2015*, pages 1491–1499.
- Ghosh, S. and Velázquez-Quesada, F. R. (2015b). A note on reliability-based preference dynamics. In van der Hoek, W., Holliday, W. H., and fan Wang, W., editors, *LORI 2015*, pages 129–142.
- Goldblatt, R. (1992). *Logics of Time and Computation*. Number 7 in CSLI Lecture Notes. Center for the Study of Language and Information, Stanford, CA, 2nd edition.
- Goldman, A. I. (2001). Experts: Which ones should you trust? *Philosophy and Phenomenological Research*, 63(1):85–110.
- Grüne-Yanoff, T. and Hansson, S. O., editors (2009). *Preference Change*, volume 42 of *Theory and Decision Library*. Springer.
- Harel, D., Kozen, D., and Tiuryn, J. (2000). *Dynamic Logic*. MIT Press, Cambridge, MA.
- Herzig, A., Lorini, E., Hübner, J. F., and Vercouter, L. (2010). A logic of trust and reputation. *Logic Journal of the IGPL*, 18(1):214–244.
- Holliday, W. H. (2010). Trust and the dynamics of testimony. In Kurzen, L., Grossi, D., and Velázquez-Quesada, F. R., editors, *Logic and Interactive Rationality. Seminar's yearbook 2009*, pages 118–142. Institute for Logic, Language and Computation, Universiteit van Amsterdam, Amsterdam, The Netherlands.
- Liau, C.-J. (2003). Belief, information acquisition, and trust in multi-agent systems – a modal logic formulation. *Artificial Intelligence*, 149(1):31–60.
- Lorini, E., Jiang, G., and Perrussel, L. (2014). Trust-based belief change. In Schaub, T., editor, *ECAI 2014 – 21st European Conference on Artificial Intelligence, 18-22 August 2014, Prague, Czech Republic – Including Prestigious Applications of Intelligent Systems (PAIS 2014)*, volume 263 of *Frontiers in Artificial Intelligence and Applications*, pages 549–554. IOS Press.
- Marx, M. and Mikuláš, S. (2001). Products, or how to create modal logics of high complexity. *Logic Journal of IGPL*, 9(1):71–82.
- Rodenhäuser, B. (2014). *A Matter of Trust: Dynamic Attitudes in Epistemic Logic*. PhD thesis, Institute for Logic, Language and Computation (ILLC), Universiteit van Amsterdam (UvA), Amsterdam, The Netherlands. ILLC Dissertation Series DS-2014-04.
- Sano, K. (2010). Axiomatizing hybrid products: How can we reason many-dimensionally in hybrid logic? *Journal of Applied Logic*, 8(4):459–474.
- Seligman, J., Liu, F., and Girard, P. (2013). Knowledge, friendship and social announcement. In van Benthem, J. and Liu, F., editors, *Logic Across the University: Foundations and Applications*, volume 47 of *Studies in Logic*, pages 445–469. College Publications.
- van Benthem, J. (2007). Dynamic logic for belief revision. *Journal of Applied Non-Classical Logics*, 17(2):129–155.
- van Benthem, J., van Eijck, J., and Kooi, B. (2006). Logics of communication and change. *Information and Computation*, 204(11):1620–1662.
- van Ditmarsch, H., van der Hoek, W., and Kooi, B. (2008). *Dynamic Epistemic Logic*. Number 337 in Synthese Library. Springer.