

# STIT based deontic logics for the miners puzzle

Xin Sun<sup>1</sup>, Zohreh Baniasadi<sup>2</sup>

<sup>1,2</sup>Faculty of Science, Technology and Communication, University of Luxembourg,  
<sup>1</sup>xin.sun@uni.lu, <sup>2</sup>zohreh.baniasadi.001@student.uni.lu

**Abstract.** In this paper we first develop two new STIT based deontic logics capable of solving the miners puzzle. The key idea is to use pessimistic lifting to lift the preference over worlds into the preference over sets of worlds. We also discuss a more general version of the miners puzzle in which plausibility is involved. In order to deal with the more general puzzle we add a modal operator representing plausibility to our logic. We present a sound and complete axiomatization.

## 1 Introduction

Research on deontic logic is divided into two main groups: the ought-to-be group and the ought-to-do group. The ought-to-do group originates from von Wright's pioneering paper [26]. Dynamic deontic logic [18, 25], deontic action logic [21, 5, 23], and STIT-based deontic logic [10, 12, 22] belong to the “ought-to-do” family.

In recent years, the miners puzzle [11] quickly grabs the attention of lots of deontic logicians [27, 6, 3, 4, 8]. The miners puzzle goes like this:

Ten miners are trapped either in shaft  $A$  or in shaft  $B$ , but we do not know which one. Water threatens to flood the shafts. We only have enough sandbags to block one shaft but not both. If one shaft is blocked, all of the water will go into the other shaft, killing every miner if they are inside. If we block neither shaft, both will be partially flooded, killing one miner.

Lacking any information about the miners' exact whereabouts, it seems acceptable to say that:

(1) We ought to block neither shaft.

However, we also accept that

(2) If the miners are in shaft  $A$ , we ought to block shaft  $A$ .

(3) If the miners are in shaft  $B$ , we ought to block shaft  $B$ .

But we also know that

(4) Either the miners are in shaft  $A$  or they are in shaft  $B$ .

And (2)-(4) seem to entail

(5) Either we ought to block shaft  $A$  or we ought to block shaft  $B$ .

Which contradicts to (1).

Various solution to this puzzle has been proposed [27, 6, 3, 4, 8]. Willer [27] claims that any adequate semantics of dyadic deontic modality must offer a solution to the miners puzzle.

The existing STIT-based deontic logic [10, 12, 22] does not offer a satisfying solution to this puzzle: although the deduction from (2)-(4) to (5) is blocked by the dyadic deontic operator defined in Sun [22], but both Horty [10] and Sun [22] are unable to predict (1). We discuss this in detail in Section 2.2.

In this paper we first develop two new STIT-based deontic logics, referring them as pessimistic utilitarian deontic logic ( $\text{PUDL}_1$  and  $\text{PUDL}_2$ ), which are capable of blocking the deduction from (2)-(4) to (5) and are able to predict (1)-(4). We further consider a more general version of the miners puzzle in which the factor of plausibility is involved. Plausibility dose not play a serious role in the original miners puzzle. It seems the plausibility of miners being in shaft  $A$  is equal to the plausibility of miner being in shaft  $B$ . If we are in a new scenario that the miners are more plausibly in shaft  $A$ , then in addition to statements (2) and (3), the following is acceptable:

- (6) We ought to block shaft  $A$ .

A logic for the miners scenario should both solve the original miners puzzle and give right predictions in the plausibility involved scenario. In this paper we extend  $\text{PUDL}_2$  to  $\text{PUDL}_2^+$  by adding a modal operator representing plausibility. We show that  $\text{PUDL}_2^+$  gives right predictions in the plausibility involved miners scenario.

The structure of this paper is as following: in Section 2 we review the existing solutions to the miners puzzle and the existing STIT-based deontic logic. In Section 3 we develop  $\text{PUDL}_1$  and  $\text{PUDL}_2$  to solve the original miners puzzle. In Section 4 we develop  $\text{PUDL}_2^+$  for the plausibility involved miners scenario. Section 5 is conclusion and future work.

## 2 Background

### 2.1 Solutions to the miners puzzle

Several authors have provided different solutions to the miners puzzle. We summarize the following approaches:

Kolodny and MacFarlane [11] give a detailed discussion of various escape routes. Then they conclude that the only possible solution to the puzzle is to invalidate the argument from (2) to (5). To do this, Kolodny and MacFarlane state we have three choices: rejecting modus ponens (MP), rejecting disjunction introduction ( $\vee I$ ), rejecting disjunction elimination ( $\vee E$ ). Among these three Kolodny and MacFarlane further demonstrate that the only wise choice is to reject MP.

Willer [27] develops a fourth option to invalidate the argument form (2) to (5): falsify the monotonicity. In his solution MP can be preserved (there are very good reasons to do so) and we are unable to derive the inconsistency.

Cariani et al [3] argue that the traditional Kratzer's semantics [13] of deontic conditionals is not capable of solving the puzzle. They propose to extend the standard

Kratzer’s account by adding a parameter representing a “decision problem” to solve the puzzle. Roughly, a decision problem contains a representation of action and a decision rule to select best action. Cariani et al [3] use a partition of all possible worlds to represent actions, and the decision rule they used to select action is essentially the same as the MaxiMin principle—the decision theoretic rule that requires agents to evaluate actions in terms of their worst conceivable outcome and choose the ‘least bad’ one among them. Such treatment shares some similarity with a special case of our logic to be in Section 3. In our logic every agent’s actions are also represented by a partition of all worlds. And we use pessimistic lifting (to be introduced later) to compare actions, which is the same as MaxiMin.

Carr [4] argues that the proposal of Cariani et al is still problematic. To develop a satisfying semantics, Carr uses three parameters to define deontic modality: an informational parameter, a value parameter and a decision rule parameter. According to Carr’s proposal, (1) to (3) are all correct predictions and no contradiction arise within her framework.

Gabbay et al [8] offers a solution to the miners puzzle using ideas from intuitionistic logic. In their logic “or” is interpreted in an intuitionistic favour. Then the deduction from statement (2), (3) and (4) to (5) is blocked.

## 2.2 STIT-based deontic logic

In STIT-based deontic logic, agents make choices and each choice is represented by a set of possible worlds. A preference relation over worlds is given as primitive. Such preference relation is then lifted to preference over sets of worlds. A choice is better than another iff the representing set of worlds of the first choice is better than the representing set of worlds of the second. A proposition  $\phi$  is obligatory (we ought to see to it that  $\phi$ ) iff it is ensured by every best choice, i.e., it is true in every world of every best choice.

Therefore the interpretation of deontic modality is based on best choices, which can only be defined on top of preference over sets of worlds. Preference over sets of worlds is defined by lifting from preference over worlds. There is no standard way of lifting preference. Lang and van der Torre [14] summarize the following three ways of lifting:

- **strong lifting:** For two sets of worlds  $W_1$  and  $W_2$ ,  $W_1$  is strongly better than  $W_2$  iff  $\forall w \in W_1, \forall v \in W_2, w$  is better than  $v$ . That is, the worst world in  $W_1$  is better than the best world in  $W_2$ .
- **optimistic lifting:**  $W_1$  is optimistically better than  $W_2$  iff  $\exists w \in W_1, \forall v \in W_2, w$  is better than  $v$ . That is, the best world in  $W_1$  is better than the best world in  $W_2$ .
- **pessimistic lifting:**  $W_1$  pessimistically better than  $W_2$  iff  $\forall w \in W_1, \exists v \in W_2, w$  is better than  $v$ . That is, the worst world in  $W_1$  is better than the worst world in  $W_2$ .

In Horty [10], Kooi and Tamminga [12] and Sun [22] the strong lifting is adopted. Applying the strong lifting to the miners scenario, all the three choices *block\_neither*, *block\_A* and *block\_B* are the best choices. “we ought to block neither” is then not true in the miners scenario. To understand this more accurately, we now give a formal review of STIT-based deontic logic of Sun [22]. We call such logic utilitarian deontic logic (UDL).

The language of the UDL is built from a finite set  $Agent = \{1, \dots, n\}$  of agents and a countable set  $\Phi = \{p, q, r, \dots\}$  of propositional letters. Let  $p \in \Phi, G \subseteq Agent$ . The UDL language  $\mathfrak{L}_{udl}$  is defined by the following Backus-Naur Form:

$$\phi ::= p \mid \neg\phi \mid \phi \wedge \phi \mid [G]\phi \mid \bigcirc_G \phi \mid \bigcirc_G(\phi/\psi)$$

Intuitively,  $[G]\phi$  is read as “group  $G$  sees to it that  $\phi$ ”.  $\bigcirc_G \phi$  is read as “ $G$  ought to see to it that  $\phi$ ”.  $\bigcirc_G(\phi/\psi)$  is read as “ $G$  ought to see to it that  $\phi$  under the condition  $\psi$ ”.

The semantics of UDL is based on utilitarian models, which is a non-temporal fragment of the group STIT model.

**Definition 1 (Utilitarian model).** A utilitarian model is a tuple  $(W, Choice, \leq, V)$ , where  $W$  is a nonempty set of possible worlds,  $Choice$  is a choice function, and  $\leq$ , representing the preference of the group  $Agent$ , is a reflexive and transitive relation over  $W$ .  $V$  is a valuation which assigns every propositional letter a set of worlds.

The choice function  $Choice : 2^{Agent} \mapsto 2^{2^W}$  is built from the individual choice function  $IndChoice : Agent \mapsto 2^{2^W}$ .  $IndChoice$  must satisfy the following conditions:

- (1) for each  $i \in Agent$  it holds that  $IndChoice(i)$  is a partition of  $W$ ;
- (2) for  $Agent = \{1, \dots, n\}$ , for every  $x_1 \in IndChoice(1), \dots, x_n \in IndChoice(n)$ ,  $x_1 \cap \dots \cap x_n \neq \emptyset$ ;

A function  $s : Agent \mapsto 2^W$  is a selection function if for each  $i \in Agent$ ,  $s(i) \in IndChoice(i)$ . Let  $Selection$  be the set of all selection functions, for every  $G \subseteq Agent$ , if  $G \neq \emptyset$ , then we define  $Choice(G) = \{\bigcap_{i \in G} s(i) : s \in Selection\}$ . If  $G = \emptyset$ , then we define  $Choice(G) = \{W\}$ .

$w \leq v$  is read as  $v$  is at least as good as  $w$ .  $w \approx v$  is short for  $w \leq v$  and  $v \leq w$ . Having defined utilitarian models, we are ready to review preferences over sets of possible worlds.

**Definition 2 (preferences over sets of worlds via strong lifting [22]).** Let  $X, Y \subseteq W$  be two sets of worlds.  $X \preceq^s Y$  ( $Y$  is at least as good as  $X$ ) if and only if

- (1) for each  $w \in X$ , for each  $w' \in Y$ ,  $w \leq w'$  and
- (2) there exists some  $v \in X$ , some  $v' \in Y$ ,  $v \leq v'$ .

$X \prec^s Y$  ( $Y$  is better than  $X$ ) if and only if  $X \preceq^s Y$  and  $Y \not\preceq^s X$ . Here the superscript  $s$  in  $\preceq^s$  is used to represent strong lifting.

**Definition 3 (dominance relation [10]).** Let  $G \subseteq Agent$  and  $K, K' \in Choice(G)$ .  $K \preceq_G^s K'$  iff for all  $S \in Choice(Agent - G)$ ,  $K \cap S \preceq^s K' \cap S$ .

$K \preceq_G^s K'$  is read as “ $K'$  weakly dominates  $K$ ”. From a decision theoretical perspective,  $K \preceq_G^s K'$  means that no matter how other agents act, the outcome of choosing  $K'$  is no worse than that of choosing  $K$ .  $K \prec_G^s K'$  is used as an abbreviation of  $K \preceq_G^s K'$  and  $K' \not\preceq_G^s K$ . If  $K \prec_G^s K'$ , then we say  $K'$  strongly dominates  $K$ .

**Definition 4 (restricted choice sets [10]).** Let  $G$  a groups of agents.

$$\text{Choice}(G/X) = \{K : K \in \text{Choice}(G) \text{ and } K \cap X \neq \emptyset\}$$

Intuitively,  $\text{Choice}(G/X)$  is the collection of those choices of group  $G$  that are consistent with condition  $X$ .

**Definition 5 (conditional dominance [22]).** Let  $G$  be a group of agents and  $X$  a set of worlds. Let  $K, K' \in \text{Choice}(G/X)$ .

$$K \preceq_{G/X}^s K' \text{ iff for all } S \in \text{Choice}((\text{Agent} - G)/(X \cap (K \cup K'))), K \cap X \cap S \preceq^s_{G/X} K' \cap X \cap S.$$

$K \preceq_{G/X}^s K'$  is read as “ $K'$  weakly dominates  $K$  under the condition of  $X$ ”. And  $K \prec_{G/X}^s K'$ , read as “ $K'$  strongly dominates  $K$  under the condition of  $X$ ”, is used to express  $K \preceq_{G/X}^s K'$  and  $K' \not\preceq_{G/X}^s K$ .

**Definition 6 (Optimal and conditional optimal [10]).** Let  $G$  be a group of agents,

- $\text{Optimal}_G^s = \{K \in \text{Choice}(G) : \text{there is no } K' \in \text{Choice}(G) \text{ such that } K \prec_G^s K'\}$ .
- $\text{Optimal}_{G/X}^s = \{K \in \text{Choice}(G/X) : \text{there's no } K' \in \text{Choice}(G/X) \text{ such that } K \prec_{G/X}^s K'\}$ .

In the semantics of UDL, the optimal choices and conditional optimal choices are used to interpret the deontic operators.

**Definition 7 (truth conditions).** Let  $M = (W, \text{choice}, \leq, V)$  be a utilitarian model and  $w \in W$ .

$$\begin{aligned} M, w \models p &\quad \text{iff } w \in V(p); \\ M, w \models \neg\phi &\quad \text{iff it is not the case that } M, w \models \phi; \\ M, w \models \phi \wedge \psi &\quad \text{iff } M, w \models \phi \text{ and } M, w \models \psi; \\ M, w \models [G]\phi &\quad \text{iff } M, w' \models \phi \text{ for all } w' \in W \text{ such that there is } K \in \text{Choice}(G), \{w, w'\} \subseteq K; \\ M, w \models \bigcirc_G \phi &\quad \text{iff } K \subseteq \|\phi\| \text{ for each } K \in \text{Optimal}_G^s; \\ M, w \models \bigcirc_G (\phi \wedge \psi) &\quad \text{iff } K \subseteq \|\phi\| \text{ for each } K \in \text{Optimal}_{G/\psi}^s. \end{aligned}$$

Here  $\|\phi\| = \{w \in W : M, w \models \phi\}$ .

**Challenge from the miners puzzle** The miners scenario is described formally by a utilitarian model as  $\text{Miners} = (W, \text{Choice}, \leq, V)$ , where  $W = \{w_1, \dots, w_6\}$ ,  $\text{Choice}(G) = \{\{w_1, w_2\}, \{w_3, w_4\}, \{w_5, w_6\}\}$ ,  $\text{Choice}(\text{Agent} - G) = \{W\}$ ,  $w_3 \approx w_6 \leq w_1 \approx w_2 \leq w_4 \approx w_5$ ,  $V(\text{in\_A}) = \{w_1, w_3, w_5\}$ ,  $V(\text{in\_B}) = \{w_2, w_4, w_6\}$ ,  $V(\text{block\_A}) = \{w_5, w_6\}$ ,  $V(\text{block\_B}) = \{w_3, w_4\}$ ,  $V(\text{block\_neither}) = \{w_1, w_2\}$ . We visualize the miners scenario by the following figure:

<i>block_neither</i>	<i>in_A</i> (9) $w_1$	$w_2$ (9) <i>in_B</i>
<i>block_B</i>	<i>in_A</i> (0) $w_3$	$w_4$ (10) <i>in_B</i>
<i>block_A</i>	<i>in_A</i> (10) $w_5$	$w_6$ (0) <i>in_B</i>

Figure 2.2:  $W = \{w_1, \dots, w_6\}$ ,  $w_3 \approx w_6 \leq w_1 \approx w_2 \leq w_4 \approx w_5$ .

Group  $G$  has three choices: *block\_neither*, *block\_A* and *block\_B*. The group of other agents has one dummy choice: choosing  $W$ . According to the semantics based on strong lifting, all the three choices are optimal. Therefore *Miners*,  $w_1 \not\models \bigcirc_G(\textit{block\_neither})$ , which means UDL fails to solve the miners puzzle.

### 3 Pessimistic utilitarian deontic logic

We now introduce pessimistic utilitarian deontic logic (PUDL) to solve the miners puzzle. We use such name because we adopt pessimistic lifting instead of strong lifting in PUDL. We develop two logics, call them PUDL<sub>1</sub> and PUDL<sub>2</sub> respectively. PUDL<sub>1</sub> is obtained from simply replacing the strong lifting in UDL by pessimistic lifting. It turns out that PUDL<sub>1</sub> is sufficient to solve the miner puzzle. But it turns out that PUDL<sub>1</sub> is bothered by other problems in deontic logic. PUDL<sub>2</sub> also solves the miners puzzle, and it is less problematic than PUDL<sub>1</sub>.

#### 3.1 PUDL<sub>1</sub>

Informally, according to the pessimistic lifting *block\_neither* is the only optimal choice in the miners scenario. Therefore “we ought to block neither” is true. It can be further proved that both (2) and (3) are true while the deduction from (2)-(4) to (5) is not valid. Therefore PUDL<sub>1</sub> offers a satisfying solution to the miners paradox. We now start to explain these arguments formally.

**Definition 8 (preferences over sets of worlds via pessimistic lifting).** Let  $X, Y \subseteq W$  be two sets of worlds.  $X \preceq^p Y$  if and only if there exists  $w \in X$ , such that for all  $w' \in Y$ ,  $w \leq w'$ .  $X \prec^p Y$  if and only if  $X \preceq^p Y$  and  $Y \not\preceq^p X$ .

**Proposition 1.**  $\preceq^p$  is reflexive and transitive.<sup>1</sup>

The pessimistic version of dominance ( $\preceq_G^p$ ), conditional dominance ( $\preceq_{G/X}^p$ ), optimal ( $Optimal_G^p$ ) and conditional optimal ( $Optimal_{G/X}^p$ ) are obtained by simply changing  $\leq^s$  to  $\leq^p$  in their strong version counterpart. We add  $\bigcirc_G^{p1}\phi$  and  $\bigcirc_G^{p1}(\phi/\psi)$  to  $\mathfrak{L}_{udl}$  to represent “from the pessimistic perspective,  $G$  ought to see to it that  $\phi$ ” and “from the pessimistic perspective,  $G$  ought to see to it that  $\phi$  in the condition  $\psi$ ” respectively. The truth condition for  $\bigcirc_G^{p1}\phi$  and  $\bigcirc_G^{p1}(\phi/\psi)$  are defined as follows:

<sup>1</sup> Due to the limitation of length, we present all proofs of propositions and theorems in the full version.

**Definition 9 (truth conditions).** Let  $M$  be a utilitarian model and  $w \in W$ .

$$\begin{aligned} M, w \models \bigcirc_G^{p_1} \phi &\quad \text{iff } K \subseteq \|\phi\| \text{ for each } K \in \text{Optimal}_G^p; \\ M, w \models \bigcirc_G^{p_1} (\phi/\psi) &\quad \text{iff } K \subseteq \|\phi\| \text{ for each } K \in \text{Optimal}_{G/\psi}^p. \end{aligned}$$

Now we return to the miners scenario. According to the pessimistic semantics, *block\_neither* is the only optimal choice. So we can draw the prediction that “we ought to block neither” i.e.  $\text{Miners}, w_1 \models \bigcirc_G^{p_1} (\text{block\_neither})$ . Moreover, given the condition of miners being in  $A$ , *block\_A* becomes the only conditional optimal choice. Hence we have “if the miners are in  $A$ , then we ought to block  $A$ ”, i.e.  $\text{Miners}, w_1 \models \bigcirc_G^{p_1} (\text{block\_A}/\text{in\_A})$ . The case for miners being in  $B$  are similar. Although we have both “if the miners are in  $A$ , then we ought to block  $A$ ” and “if the miners are in  $B$ , then we ought to block  $B$ ”, by Proposition 2 below we can avoid the prediction that “we ought to block either  $A$  or  $B$ ”. Hence no contradiction arise. Therefore PUDL<sub>1</sub> gives right prediction meanwhile avoids contradictions. It therefore offers a viable solution to the miners puzzle.

**Proposition 2.**  $\nexists \bigcirc_G^{p_1} (p/q) \wedge \bigcirc_G^{p_1} (p/r) \rightarrow \bigcirc_G^{p_1} (p/(q \vee r))$ .

### 3.2 PUDL<sub>2</sub>

Although PUDL<sub>1</sub> solves the miners puzzle, it still has some drawbacks. On the intuitive side, PUDL<sub>1</sub> is not free from Ross’ paradox. Ross’ paradox [19] originate from the logic of imperatives, and is a well-known puzzle in deontic logic which can be concisely stated as following:

Suppose you ought to mail the letter. Since mail the letter logically entails mail the letter or burn it, you ought to mail the letter or burn it.

PUDL<sub>1</sub> validates the formula  $\bigcirc_G^{p_1} p \rightarrow \bigcirc_G^{p_1} (p \vee q)$ , which means it is not free from Ross’ paradox.

On the technical side, PUDL<sub>1</sub> is not finitely axiomatizable. This is because PUDL<sub>1</sub> contains group STIT. Herzig and Schwarzenbuber [9] show that if  $|\text{Agent}| \geq 3$  then group STIT is not finitely axiomatizable.

To fix these flaws, we develop PUDL<sub>2</sub>. We show that PUDL<sub>2</sub> solves the miners puzzle and is free from the Ross’s paradox. We further give an axiomatization of PUDL<sub>2</sub>.

**Language** Similar to  $\mathcal{L}_{udl}$ , the language of the PUDL<sub>2</sub> is built from  $\text{Agent}$  and  $\Phi$ . But for the sake of axiomatization, we simplify group STIT in UDL to individual STIT. In order to syntactically define pessimistic lifting we add a preference modality as well as the universal modality to our language. For  $p, q \in \Phi$  and  $i \in \text{Agent}$ , the language  $\mathcal{L}_{pudl}^2$  is given by the following Backus-Naur Form:

$$\phi ::= p \mid \neg\phi \mid \phi \wedge \phi \mid [i]\phi \mid \Box\phi \mid [\leq]\phi \mid [\geq]\phi \mid [ < ]\phi$$

Intuitively,  $[i]\phi$  means “agent  $i$  sees to it that  $\phi$ ”.  $\Box\phi$  means “ $\phi$  is true everywhere”.  $[\leq]\phi$  means “ $\phi$  is weakly preferable” while  $[ < ]\phi$  means “ $\phi$  is strictly preferable”.  $[\geq]\phi$

means “ $\phi$  is unpreferable”. We use  $\diamond$ ,  $\langle \leq \rangle$  and  $\langle < \rangle$  as the dual for  $\square$ ,  $\leq$  and  $<$  respectively.

Semantically the preference relation  $\leq$  corresponding to  $\leq$  is required to be a weak linear order. That is,  $\leq$  is reflexive, transitive and connected. The preference relation  $<$  corresponding to  $<$  is required to satisfy the following:  $w < v$  iff  $w \leq v$  and  $v \not\leq w$ . Lifting of preference can be defined in  $\mathcal{L}_{pudl}^2$  only with these constraints. Liu [16] observes that it is sufficient to define optimistic lifting with  $\leq$  being partial order. But to define strong and pessimistic lifting,  $\leq$  is required to be linear.

- strong lifting:  $\phi \leq^s \psi ::= \square(\psi \rightarrow [ < ] \neg \phi)$ . Intuitively,  $\square(\psi \rightarrow [ < ] \neg \phi)$  says that for all  $\psi$ -world, there is no  $\phi$  world which is better. In other words, every  $\psi$ -world is at least as good as every  $\phi$ -world. That is, the worst  $\psi$ -world is at least as good as the best  $\phi$ -world.
- optimistic lifting:  $\phi \leq^o \psi ::= \square(\phi \rightarrow \langle \leq \rangle \psi)$ . Intuitively,  $\square(\phi \rightarrow \langle \leq \rangle \psi)$  says that for all  $\phi$ -world  $w$  there is a  $\psi$ -world which is at least as good as  $w$ . In other words, for the best  $\phi$ -world  $w$  there is a  $\psi$ -world which is at least as good as  $w$ . That is, the best  $\psi$ -world is at least as good as the best  $\phi$ -world.
- pessimistic lifting:  $\phi \leq^p \psi ::= \square(\psi \rightarrow \langle \geq \rangle \phi)$ . Intuitively,  $\square(\psi \rightarrow \langle \geq \rangle \phi)$  says that for all  $\psi$ -world  $w$ , it is at least as good as some  $\phi$ -world. That is, the worst  $\psi$ -world is at least as good as the worst  $\phi$ -world.<sup>2</sup>

We use  $\phi <^p \psi$  as an abbreviation of  $(\phi \leq^p \psi) \wedge \neg(\psi \leq^p \phi)$ . Obligation and conditional obligation are defined in our language as follows:

- $\bigcirc_i^{p2} \phi ::= \diamond[i]\phi \wedge (\neg\phi <^p [i]\phi)$ . Intuitively, agent  $i$  is obligatory to see to it that  $\phi$  iff it is possible for  $i$  to see to it that  $\phi$  and seeing to it that  $\phi$  is strictly better than  $\neg\phi$  in the pessimistic sense.
- $\bigcirc_i^{p2}(\phi/\psi) ::= \diamond[i]\phi \wedge ((\neg\phi \wedge \psi) <^p ([i]\phi \wedge \psi))$ .

**Semantics** The semantics of pessimistic utilitarian deontic logic is based on pessimistic utilitarian model, which is a non-temporal individual fragment of the STIT model.

**Definition 10 (Pessimistic utilitarian model).** A pessimistic utilitarian model is a tuple  $M = (W, IndChoice, \leq, <, V)$ , where  $W$  is a nonempty set of possible worlds,  $IndChoice$  is an individual choice function,  $\leq$  is a reflexive, transitive and connected relation over  $W$ , representing the preference of the group Agent.  $<$  is a sub-relation of  $\leq$  such that for all  $w, w' \in W$ ,  $w < w'$  iff  $w \leq w'$  and  $w' \not\leq w$ .

The individual choice function  $IndChoice : Agent \mapsto 2^{2^W}$  must satisfy the following conditions:

- (1) for each  $i \in Agent$  it holds that  $IndChoice(i)$  is a partition of  $W$ ;
- (2) for  $Agent = \{1, \dots, n\}$ , for every  $x_1 \in IndChoice(1), \dots, x_n \in IndChoice(n)$ ,  $x_1 \cap \dots \cap x_n \neq \emptyset$ ;

---

<sup>2</sup> In Definition 3.8 of Liu [16], pessimistic lifting  $\phi \leq^p \psi$  is defined as  $\diamond(\phi \wedge [ < ] \neg \psi)$ . Such treatment is problematic because we can construct a counter example such that  $\phi \leq^p \psi$  is true but the worst  $\psi$ -world is NOT at least as good as the worst  $\phi$ -world. Here is a counter example:  $W = \{w_1, w_2\}$ ,  $w_1 < w_2$ ,  $\phi = \{w_2\}$ ,  $\|\psi\| = \{w_1, w_2\}$ .

Let  $R_i$  be the equivalence relation induced by  $IndChoice(i)$ . Then  $(w, w') \in R_i$  iff there is  $K \in IndChoice(i)$  such that  $\{w, w'\} \subseteq K$ .  $IndChoice(i) = \{R_i(w) : w \in W\}$ , where  $R_i(w) = \{w' \in W : (w, w') \in R_i\}$ . The truth condition of formulas of  $\mathcal{L}_{pudl}^2$  is defied as follows:

**Definition 11 (truth conditions).** Let  $M$  be a pessimistic utilitarian model,  $w \in W$ .

- $M, w \models_{pudl_2} [i]\phi$  iff  $M, w' \models \phi$  for all  $w'$  such that  $(w, w') \in R_i$ ;
- $M, w \models_{pudl_2} [\leq]\phi$  iff  $M, w' \models \phi$  for all  $w'$  such that  $w \leq w'$ ;
- $M, w \models_{pudl_2} [\geq]\phi$  iff  $M, w' \models \phi$  for all  $w'$  such that  $w' \leq w$ ;
- $M, w \models_{pudl_2} [<]\phi$  iff  $M, w' \models \phi$  for all  $w'$  such that  $w < w'$ ;
- $M, w \models_{pudl_2} \Box\phi$  iff  $M, w' \models \phi$  for all  $w' \in W$ .

The axiomatization of  $PUDL_2$  is a fragment of the axiomatization of  $PUDL^+$  in the next section. The following proposition shows that  $PUDL_2$  is free from Ross' paradox.

**Proposition 3.**  $\nexists_{pudl_2} \bigcirc_i^{p_2} p \rightarrow \bigcirc_i^{p_2} (p \vee q)$ .

**Another analysis to the miners puzzle** The miners scenario is described formally by a pessimistic utilitarian model as  $Miners^p = (W, IndChoice, \leq, <, V)$ , where  $W = \{w_1, \dots, w_6\}$ ,  $IndChoice(i) = \{\{w_1, w_2\}, \{w_3, w_4\}, \{w_5, w_6\}\}$ ,  $IndChoice(j) = \{W\}$  for all  $j \neq i$ ,  $w_3 \approx w_6 < w_1 \approx w_2 < w_4 \approx w_5$ ,  $V(in\_A) = \{w_1, w_3, w_5\}$ ,  $V(in\_B) = \{w_2, w_4, w_6\}$ ,  $V(block\_A) = \{w_5, w_6\}$ ,  $V(block\_B) = \{w_3, w_4\}$ ,  $V(block\_neither) = \{w_1, w_2\}$ .

Agent  $i$  is able to see to it that:  $block\_neither$ ,  $block\_A$  and  $block\_B$ .  $[i]block\_neither$  is true in worlds  $w_1$  and  $w_2$ . According to the pessimistic semantics,  $[i]block\_neither$  is strictly better than  $\neg block\_neither$ . Therefore  $i$  ought to block neither. That is,  $Miners^p, w_1 \models \bigcirc_G^{p_2} (block\_neither)$ .

Moreover, given the condition of miners being in  $A$ ,  $[i]block\_A$  is better than  $\neg block\_A$ . Hence we have “if the miners are in  $A$ , then  $i$  ought to block  $A$ ”. That is,  $Miners^p, w_1 \models \bigcirc_i^{p_2} (block\_A/in\_A)$ . The case for miners being in  $B$  is similar.

It remains to show that although we have both “if the miners are in  $A$ , then we ought to block  $A$ ” and “if the miners are in  $B$ , then we ought to block  $B$ ”, but we cannot logically derive “we ought to block either  $A$  or  $B$ ”. This is done by the following proposition.

**Proposition 4.**  $\nexists_{pudl_2} \bigcirc_i^{p_2} (p/q) \wedge \bigcirc_i^{p_2} (p/r) \rightarrow \bigcirc_i^{p_2} (p/(q \vee r))$

## 4 Plausiblity involved pessimistic utilitarian deontic logic

The interplay of plausibility and preference are heavily discussed in qualitative decision theory [2, 7]. Boutilier [1] uses the modality of plausibility and preference to define conditional goals. Lang *et al* [15] use plausibility and preference to define hidden desire.

In this section we develop plausiblity involved pessimistic utilitarian deontic logic  $PUDL_2^+$  to analyze the plausibility involved miners puzzle. The language of  $PUDL_2^+$  is  $\mathcal{L}_{pudl}^2$  extended with plausibility operators. Formally, for  $p, q \in \Phi$  and  $i \in Agent$ , the language  $\mathcal{L}_{pudl}^{2+}$  is given by the following Backus-Naur Form:

$$\phi ::= p \mid \neg\phi \mid \phi \wedge \phi \mid [i]\phi \mid \Box\phi \mid [\leq]\phi \mid [\geq]\phi \mid [ < ]\phi \mid [\leq_p]\phi \mid [ <_p ]\phi$$

Plausibility involved pessimistic lifting is defined as follows:

$$\phi \leq_p^p \psi ::= (\phi \wedge [ <_p ]\neg\phi) \leq^p (\psi \wedge [ <_p ]\neg\psi)$$

Intuitively,  $\phi \leq_p^p \psi$  says that the most plausible  $\psi$  worlds are better than the most plausible  $\phi$  worlds from a pessimistic perspective. We use  $\phi <_p^p \psi$  as an abbreviation of  $(\phi \leq_p^p \psi) \wedge \neg(\psi \leq_p^p \phi)$ . Plausibility involved obligation and conditional obligation are defined in  $\mathfrak{L}_{pudl}^{2+}$  as follows:

- $\bigcirc_i \phi ::= \Diamond[i]\phi \wedge (\neg\phi <_p^p [i]\phi)$ .
- $\bigcirc_i(\phi/\psi) ::= \Diamond[i]\phi \wedge ((\neg\phi \wedge \psi) <_p^p ([i]\phi \wedge \psi))$ .

#### 4.1 Semantics

**Definition 12 (Plausibility involved pessimistic utilitarian model).** A *plausibility involved pessimistic utilitarian model* is a tuple  $(W, IndChoice, \leq, <, \leq_p, <_p, V)$ , where  $(W, IndChoice, \leq, <, V)$  is a pessimistic utilitarian model.  $\leq_p$  is a reflexive, transitive and connected relation over  $W$ , representing plausibility.  $<_p$  is a sub-relation of  $\leq_p$  such that for all  $w, w' \in W$ ,  $w <_p w'$  iff  $w \leq_p w'$  and  $w' \not\leq_p w$ .

The truth condition of formulas in  $\mathfrak{L}_{pudl}^{2+}$  is the same as  $\mathfrak{L}_{pudl}^2$ , except those formulas contains plausibility operators.

**Definition 13 (truth conditions).** Let  $M$  be a pessimistic utilitarian model,  $w \in W$ .

$$\begin{aligned} M, w \models_{pudl_2^+} [\leq_p]\phi &\text{ iff } M, w' \models \phi \text{ for all } w' \text{ such that } w \leq_p w'; \\ M, w \models_{pudl_2^+} [ <_p ]\phi &\text{ iff } M, w' \models \phi \text{ for all } w' \text{ such that } w <_p w'; \end{aligned}$$

In the generalized miners puzzle. Since it is more plausible that miners are in shaft  $A$ ,  $block\_A$  is the only optimal choice. Therefore  $\bigcirc_i block\_A$  is true. Given the miner are in  $B$ ,  $block\_B$  is the conditional optimal choice, therefore  $\bigcirc_i(block\_B/in\_B)$ .

$block\_neither$	$in\_A(9)$	$w_1$	$w_2$	$(9)$	$in\_B$
	$in\_A(0)$	$w_3$	$w_4$	$(10)$	$in\_B$
	$in\_A(10)$	$w_5$	$w_6$	$(0)$	$in\_B$

Figure 4.1:  $W = \{w_1, \dots, w_6\}$ ,  $w_3 \approx w_6 \leq w_1 \approx w_2 \leq w_4 \approx w_5$ ,  
 $w_2 \approx_p w_4 \approx_p w_6 <_p w_1 \approx_p w_3 \approx_p w_5$ .

## 4.2 Proof system

The proof system of  $\text{PUDL}_2^+$  consists the following axioms and the rules of *modus ponens*, and *necessitation* for  $[1], \dots, [n], \square, [\leq], [\geq], [ < ], [\leq_p]$  and  $[ <_p ]$ . The following is the list of axioms:

1. S4.3 for  $[\leq]$ 
  - (a)  $[\leq](\phi \rightarrow \psi) \rightarrow ([\leq]\phi \rightarrow [\leq]\psi)$
  - (b)  $[\leq]\phi \rightarrow \phi$
  - (c)  $[\leq]\phi \rightarrow [\leq][\leq]\phi$
  - (d)  $\langle \leq \rangle \phi \wedge \langle \leq \rangle \psi \rightarrow (\langle \leq \rangle (\phi \wedge \langle \leq \rangle \psi) \vee \langle \leq \rangle (\phi \wedge \psi) \vee \langle \leq \rangle (\psi \wedge \langle \leq \rangle \phi))$
2. S4.3 for  $[\leq_p]$ 
  - (a)  $[\leq_p](\phi \rightarrow \psi) \rightarrow ([\leq_p]\phi \rightarrow [\leq_p]\psi)$
  - (b)  $[\leq_p]\phi \rightarrow \phi$
  - (c)  $[\leq_p]\phi \rightarrow [\leq_p][\leq_p]\phi$
  - (d)  $\langle \leq_p \rangle \phi \wedge \langle \leq_p \rangle \psi \rightarrow (\langle \leq_p \rangle (\phi \wedge \langle \leq_p \rangle \psi) \vee \langle \leq_p \rangle (\phi \wedge \psi) \vee \langle \leq_p \rangle (\psi \wedge \langle \leq_p \rangle \phi))$
3. Mutual converse for  $[\leq]$  and  $[\geq]$ :
$$(\phi \rightarrow [\leq]\langle \geq \rangle \phi) \wedge (\phi \rightarrow [\geq]\langle \leq \rangle \phi)$$
4. K for  $[ < ]$ :
$$[ < ](\phi \rightarrow \psi) \rightarrow ([ < ]\phi \rightarrow [ < ]\psi)$$
5. K for  $[ <_p ]$ :
$$[ <_p ](\phi \rightarrow \psi) \rightarrow ([ <_p ]\phi \rightarrow [ <_p ]\psi)$$
6. Interaction
  - (a)  $[ < ]\phi \rightarrow [ < ][\leq]\phi$
  - (b)  $[ < ]\phi \rightarrow [\leq][ < ]\phi$
  - (c)  $[\leq]([\leq]\phi \vee \psi) \wedge [ < ]\psi \rightarrow \phi \vee [\leq]\psi$
  - (d)  $[ <_p ]\phi \rightarrow [ <_p ][\leq_p]\phi$
  - (e)  $[ <_p ]\phi \rightarrow [\leq_p][ <_p ]\phi$
  - (f)  $[\leq_p]([\leq_p]\phi \vee \psi) \wedge [ <_p ]\psi \rightarrow \phi \vee [\leq_p]\psi$
7. Inclusion
  - (a)  $[\leq]\phi \rightarrow [ < ]\phi$
  - (b)  $\square\phi \rightarrow [\leq]\phi$
  - (c)  $[\leq_p]\phi \rightarrow [ <_p ]\phi$
  - (d)  $\square\phi \rightarrow [\leq_p]\phi$
  - (e)  $\square\phi \rightarrow [i]\phi$ , for  $i \in \text{Agent}$
8. S5 for  $\square$  and  $[i]$ ,  $i \in \text{Agent}$
9. Agent independent:  $(\Diamond[1]\phi_1 \wedge \dots \wedge \Diamond[n]\phi_n) \rightarrow \Diamond([1]\phi_1 \wedge \dots \wedge [n]\phi_n)$

For every  $\phi$  is derivable from the proof system of  $\text{PUDL}_2^+$ , then we say  $\phi$  is a theorem of  $\text{PUDL}_2^+$  and write  $\vdash \phi$ . For a set of formulas  $\Gamma \cup \phi$ , we say  $\phi$  is derivable from  $\Gamma$  (write  $\Gamma \vdash \phi$ ) if  $\vdash \phi$  or there are formulas  $\psi_1, \dots, \psi_n \in \Gamma$  such that  $\vdash (\psi_1 \wedge \dots \wedge \psi_n) \rightarrow \phi$ .

**Theorem 1 (soundness and completeness).**  $\Gamma \vdash \phi$  iff  $\Gamma \models_{\text{pudl}_2^+} \phi$

The proof of soundness is routine. For completeness, we adopt the canonical model method in addition with Bulldozing [20]: we first build a canonical model, then we transform the canonical model via Bulldozing to make a new model to satisfy the requirement of plausibility involved pessimistic utilitarian model. We sketch the proof in the appendix.

## 5 Conclusion and future work

In this paper we first develop two new STIT based deontic logics capable of solving the miners puzzle. The key idea is to use pessimistic lifting to lift preference over worlds to preference over sets of worlds. To deal with the more general miners scenario we add modal operators representing plausibility. A complete axiomatization is given. Concerning future works, the most natural extension is to replace non-temporal STIT by temporal STIT logic [17].

## References

1. Craig Boutilier. Toward a logic for qualitative decision theory. In Jon Doyle, Erik Sandewall, and Pietro Torasso, editors, *Proceedings of the 4th International Conference on Principles of Knowledge Representation and Reasoning*, pages 75–86, Bonn, Germany, 1994. Morgan Kaufmann.
2. Ronen Brafman and Moshe Tennenholtz. Modeling agents as qualitative decision makers. *Artificial Intelligence*, 94(1-2):217–268, 1997.
3. Fabrizio Cariani, Magdalena Kaufmann, and Stefan Kaufmann. Deliberative modality under epistemic uncertainty. *Linguistics and Philosophy*, 36(3):225–259, 2013.
4. Jennifer Carr. Deontic modals without decision theory. *Proceedings of Sinn und Bedeutung*, 17:167–182, 2012.
5. Pablo Castro and T.S.E. Maibaum. Deontic action logic, atomic boolean algebras and fault-tolerance. *Journal of Applied Logic*, 2009.
6. Nate Charlow. What we know and what to do. *Synthese*, 190(12):2291–2323, 2013.
7. Jon Doyle and Richmond Thomason. Background to qualitative decision theory. *AI Magazine*, 20(2):55–68, 1999.
8. Dov M. Gabbay, Livio Robaldo, Xin Sun, Leendert van der Torre, and Zohreh Baniasadi. Toward a linguistic interpretation of deontic paradoxes - beth-reichenbach semantics approach for a new analysis of the miners scenario. In *Deontic Logic and Normative Systems - 12th International Conference, DEON 2014, Ghent, Belgium, July 12-15, 2014. Proceedings*, pages 108–123, 2014.
9. Andreas Herzig and François Schwarzentruber. Properties of logics of individual and group agency. In Carlos Areces and Robert Goldblatt, editors, *Advances in Modal Logic 7, papers from the seventh conference on "Advances in Modal Logic," held in Nancy, France, 9-12 September 2008*, pages 133–149. College Publications, 2008.
10. John Horty. *Agency and Deontic Logic*. Oxford University Press, New York, 2001.
11. Niko Kolodny and John MacFarlane. Iffs and oughts. *Journal of Philosophy*, 107(3):115–143, 2010.
12. Barteld Kooi and Allard Tamminga. Moral conflicts between groups of agents. *Journal of Philosophical Logic*, 37:1–21, 2008.
13. Angelika Kratzer. The notional category of modality. In H. J. Eikmeyer and H. Rieser, editors, *Words, worlds, and Contexts: New Approaches in World Semantics*. Berlin: de Gruyter, 1981.
14. Jerom Lang and Leendert van der Torre. From belief change to preference change. In M. Ghallab, C.D. Spyropoulos, N. Fakotakis, and N. Avouris, editors, *Proceedings of the 2008 conference on ECAI 2008: 18th European Conference on Artificial Intelligence*, pages 351–355, Amsterdam, 2008. IOS Press.
15. Jerom Lang, Leendert van der Torre, and Emil Weydert. Hidden uncertainty in the logical representation of desires. *Proceedings of IJCAI2003*, pages 685–690, 2003.

16. Fenrong Liu. *Reasoning about Preference Dynamics*. Springer, 2011.
17. Emiliano Lorini. Temporal stit logic and its application to normative reasoning. *Journal of Applied Non-Classical Logics*, 23(4):372–399, 2013.
18. John Jule Meyer. A different approach to deontic logic: deontic logic viewed as a variant of dynamic logic. *Notre Dame Journal of Formal Logic*, pages 109–136, 1988.
19. A. Ross. Imperatives and logic. *Theoria*, 7(5371), 1941.
20. Krister Segerberg. *An essay in classical modal logic*, volume 13 of *Filosofiska Studier*. Uppsala: Filosofiska foreningen och Filosofiska institutionen vid Uppsala universitet, 1971.
21. Krister Segerberg. A deontic logic of action. *Studia Logica*, 1982.
22. Xin Sun. Conditional ought, a game theoretical perspective. In J. Lang H. van Ditmarsch and S. Ju, editors, *Logic, Rationality, and Interaction: Proceedings of the Thire International Workshop*, pages 356–369, Guangzhou, China, 2011.
23. Robert Trypuz and Piotr Kulicki. On deontic action logics based on boolean algebra. *Journal of Logic and Computation*, [forthcoming], 2014.
24. Johan van Benthem, Patric Girard, and Olivier Roy. Evernthing else being equal: A modal logic approach for ceteris paribus preference. *Journal of Philosophical Logic*, 38(1):83–125, 2009.
25. Ron van der Meyden. The dynamic logic of permission. *Journal of Logic and Computation*, 6:465–479, 1996.
26. Georg von Wright. Deontic logic. *Mind*, pages 1–15, 1951.
27. Malte Willer. A remark on iffy oughts. *Journal of Philosophy*, 109(7):449461, 2012.

## Appendix

**Proposition 5.** *If every consistent  $\Gamma$  is satisfiable on some model  $M$ , then  $\Gamma \models_{pudi_2^+} \phi$  implies  $\Gamma \vdash \phi$ .*

**Definition 14 (maximal consistent set (MCS)).** *A set of formulas  $\Gamma$  is maximal consistent if  $\Gamma$  is consistent and any proper extension of  $\Gamma$  is not consistent.*

For every consistent  $\Gamma$ ,  $\Gamma$  can be extend to a MCS  $\Gamma^+$ , we then construct a canonical model for  $\Gamma^+$

**Definition 15 (canonical model).** *The canonical model  $\mathfrak{M}^0$  for  $\Gamma^+$  is a relational structure  $(W^0, \{R_i^0\}_{i \in \text{Agent}}, \leq^0, \leq_p^0, <^0, V^0)$  where:*

- $W^0 = \{w \mid w \text{ is a MCS and for all } \Box\phi \in \Gamma^+, \phi \in w\}$ ;
- For every  $i \in \text{Agent}$ ,  $R_i^0$  is a binary relation on  $W^0$  defined by  $w R_i^0 v$  iff for all  $\phi$ ,  $[i]\phi \in w$  implies  $\phi \in v$ ;
- $\leq^0$  is a binary relation on  $W^0$  defined by  $w \leq^0 v$  iff for all  $\phi$ ,  $[\leq]\phi \in w$  implies  $\phi \in v$ ;
- $<^0$  is a binary relation on  $W^0$  defined by  $w <^0 v$  iff for all  $\phi$ ,  $[<]\phi \in w$  implies  $\phi \in v$ ;
- $\leq_p^0$  is a binary relation on  $W^0$  defined by  $w \leq_p^0 v$  iff for all  $\phi$ ,  $[\leq_p]\phi \in w$  implies  $\phi \in v$ ;
- $<_p^0$  is a binary relation on  $W^0$  defined by  $w <_p^0 v$  iff for all  $\phi$ ,  $[<_p]\phi \in w$  implies  $\phi \in v$ ;
- $V^0$  is the valuation defined by  $V^0(p) = \{w \in W^0 \mid p \in w\}$ .

**Proposition 6.**  $\mathfrak{M}^0, \Gamma^+ \models_{pudl_2^+} \Gamma$ .

**Proposition 7.**  $\mathfrak{M}^0$  has the following properties:

- (1) Both  $\leq^0$  and  $\leq_p^0$  are reflexive, transitive and connected relations.
- (2) If  $w \leq^0 v$  and  $v \not\leq^0 w$  then  $w <^0 v$ .
- (3) If  $w \leq_p^0 v$  and  $v \not\leq_p^0 w$  then  $w <_p^0 v$ .
- (4) If  $w <^0 v$  then  $w \leq^0 v$ .
- (5) If  $w <_p^0 v$  then  $w \leq_p^0 v$ .
- (6)  $R_i^0$  is an equivalence relation for each  $i \in \text{Agent}$ .
- (7) For every  $w \in W^0$ ,  $R_1^0(w) \cap \dots \cap R_n^0(w) \neq \emptyset$ .

**Deleting  $<$ -cluster** Note that converse of item (2) of Proposition 7 is not true because there may be two world  $w$  and  $v$  in  $W^0$  such that  $w <^0 v$  and  $v <^0 w$ . In this case we say that  $w$  and  $v$  are in the same  $<^0$ -clusters. To deal with this we follow Bentham [24] to use the technique called Bulldozing [20] to transform  $\mathfrak{M}^0$  to a new model  $\mathfrak{M}^1$  such that there is no  $<$ -cluster in  $\mathfrak{M}^1$ .

**Definition 16 (cluster).** A  $<$ -cluster is an inclusion-maximal set of worlds  $C$  such that  $w < v$  for all worlds  $w, v \in C$ . Similarly for  $\leq_p$ -cluster.

Let  $\mathfrak{M}^1 = (W^1, \{R_i^1\}_{i \in \text{Agent}}, \leq^1, <^1, \leq_p^1, <_p^1, V^1)$  where:

- $W^1 = W^{0-} \cup \bigcup_{i \in I} C'_i$ , here  $I$  is a set index of all  $<$ -clusters of  $W^0$ ,  $W^{0-} = W^0 - \bigcup_{i \in I} C_i$ ,  $C'_i = C_i \times \mathbb{Z}$ ,  $\mathbb{Z}$  is the set of natural numbers.
- $R_i^1$  is defined by  $w R_i^1 v$  iff  $\beta(w) R_i^0 \beta(v)$ , for every  $i \in \text{Agent}$ .
- $<^1$  is defined as follows: For each  $C_i$ , choose an arbitrary linear order  $<^{1,i}$ . Define a map  $\beta : W^1 \rightarrow W^0$  by  $\beta(x) = x$  if  $x \in W^{0-}$  and  $\beta(x) = w$  if  $x$  is a pair  $(w, n)$  for some world  $w$  and integer  $n$ . We define  $<^1$  via the following cases:
  - Case 1:  $x$  or  $y$  is in  $W^{0-}$ . In this case we let  $x <^1 y$  iff  $\beta(x) <^0 \beta(y)$ .
  - Case 2:  $x \in C'_i$  and  $y \in C'_j$ ,  $i \neq j$ . In this case we let  $x <^1 y$  iff  $\beta(x) <^0 \beta(y)$ .
  - Case 3:  $x \in C'_i$  and  $y \in C'_i$ . In this case,  $x = (w, m)$  and  $y = (v, n)$ . There are two sub-cases:
    - \* Case 3.1: If  $m \neq n$ , we use the natural ordering on  $\mathbb{Z}$ :  $(w, m) <^1 (v, n)$  iff  $m < n$ .
    - \* Case 3.2: If  $m = n$ , we use the linear ordering  $<^{1,i}$ :  $(w, m) <^1 (v, m)$  iff  $w <^{1,i} v$ .
- $\leq^1$  is defined via the following cases:
  - Case 1:  $x$  or  $y$  is in  $W^{0-}$ . In this case we let  $x \leq^1 y$  iff  $\beta(x) \leq^0 \beta(y)$ .
  - Case 2: Otherwise, we take the reflexive closure of  $<^1$ :  $x \leq^1 y$  iff  $x <^1 y$  or  $x = y$ .
- $\leq_p^1$  is defined by  $w \leq_p^1 v$  iff  $\beta(w) \leq_p^0 \beta(v)$ .
- $<_p^1$  is defined by  $w <_p^1 v$  iff  $\beta(w) <_p^0 \beta(v)$ .
- $V^1$  is defined by  $w \in V^1(p)$  iff  $\beta(w) \in V^0(p)$ .

**Definition 17 (Bounded Morphism).** A mapping  $f : M = (W, \{R_i\}_{i \in \text{Agent}}, \leq, <, \leq_p, <_p, V) \rightarrow M' = (W, \{R'_i\}_{i \in \text{Agent}}, \leq', <', \leq'_p, <'_p, V')$  is a bounded morphism if it satisfies the following conditions:

- $w$  and  $f(w)$  satisfy the same proposition letters.
- if  $w \leq v$  then  $f(w) \leq' f(v)$ . And similarly for  $<, \leq_p, \leq'_p, R_i$ .
- if  $f(w) \leq' v'$  then there exists  $v$  such that  $w \leq v$  and  $f(v) = v'$ . And similarly for  $<', \leq'_p, \leq'_p, R_i$ .

**Lemma 1.** If  $f$  is a bounded morphism from  $M$  to  $M'$ , then for all  $\phi$ , for all  $w \in M$ ,  $M, w \models_{pudl_2^+} \phi$  iff  $M', f(w) \models_{pudl_2^+} \phi$ .

**Proposition 8.** For every consistent set  $\Phi$ , if  $\mathfrak{M}^0, \Gamma \models_{pudl_2^+} \Phi$ , then there exist  $\Gamma'$  such that  $\mathfrak{M}^1, \Gamma' \models_{pudl_2^+} \Phi$ .

**Proposition 9.**  $\mathfrak{M}^1$  has the following properties:

- (1) Both  $\leq^1$  and  $\leq_p^1$  are reflexive, transitive and connected relations.
- (2)  $w <^1 v$  iff  $w \leq^1 v$  and  $v \not\leq^1 w$
- (3) If  $w \leq_p^1 v$  and  $v \not\leq_p^1 w$  then  $w <_p^1 v$ .
- (4) If  $w <_p^1 v$  then  $w \leq_p^1 v$ .
- (5)  $R_i^1$  is an equivalence relation for each  $i \in \text{Agent}$ .
- (6) For every  $w \in W^1$ ,  $R_1^1(w) \cap \dots \cap R_n^1(w) \neq \emptyset$ .

**Deleting  $\leq_p$ -cluster** Now we use bulldozing again to delete  $\leq_G$  clusters.

Let  $\mathfrak{M}^2 = (W^2, \{R_i^2\}_{i \in \text{Agent}}, \leq^2, <^2, \leq_p^2, V^2)$  where:

- $W^2 = W^{1-} \cup \bigcup_{i \in I} C'_i$ , here  $I$  is a set index of all  $\leq_G$ -clusters of  $W^1$ ,  $W^{1-} = W^1 - \bigcup_{i \in I} C_i$ ,  $C'_i = C_i \times \mathbb{Z}$ ,  $\mathbb{Z}$  is the set of natural numbers.
- $R_i^2$  is defined by  $w R_i^2 v$  iff  $\sigma(w) R_i^1 \sigma(v)$ , for every  $i \in \text{Agent}$ .
- $<_p^2$  is defined as follows: For each  $C_i$ , choose an arbitrary linear order  $<_p^{2,i}$ . Define a map  $\sigma : W^2 \rightarrow W^1$  by  $\sigma(x) = x$  if  $x \in W^{1-}$  and  $\sigma(x) = w$  if  $x$  is a pair  $(w, n)$  for some world  $w$  and integer  $n$ . We define  $<_p^2$  via the following cases:
  - Case 1:  $x$  or  $y$  is in  $W^{1-}$ . In this case we let  $x <_p^2 y$  iff  $\sigma(x) <_p^1 \sigma(y)$ .
  - Case 2:  $x \in C_i$  and  $y \in C_j$ ,  $i \neq j$ . In this case we let  $x <_p^3 y$  iff  $\sigma(x) <_p^2 \sigma(y)$ .
  - Case 3:  $x \in C_i$  and  $y \in C_i$ . In this case,  $x = (w, m)$  and  $y = (v, n)$ . There are two sub-cases:
    - \* Case 3.1: If  $m \neq n$ , we use the natural ordering on  $\mathbb{Z}$ :  $(w, m) <_p^2 (v, n)$  iff  $m < n$ .
    - \* Case 3.2: If  $m = n$ , we use the linear ordering  $<_p^{2,i}$ :  $(w, m) <_p^2 (v, m)$  iff  $w <_p^{2,i} v$ .
- $\leq_p^2$  is defined via the following cases:
  - Case 1:  $x$  or  $y$  is in  $W^{1-}$ . In this case we let  $x \leq_p^2 y$  iff  $\sigma(x) \leq_p^1 \sigma(y)$ .
  - Case 2: Otherwise, we take the reflexive closure of  $<_p^2$ :  $x \leq_p^2 y$  iff  $x <_p^2 y$  or  $x = y$ .
- $\leq^2$  is defined by  $w \leq^2 v$  iff  $\sigma(w) \leq^1 \sigma(v)$ .
- $<^2$  is defined by  $w <^2 v$  iff  $\sigma(w) <^1 \sigma(v)$ .
- $V^2$  is defined by  $w \in V^2(p)$  iff  $\sigma(w) \in V^1(p)$ .

**Proposition 10.** For every consistent set  $\Phi$ , if  $\mathfrak{M}^1, \Gamma \models_{pudl_2^+} \Phi$ , then there exist  $\Gamma'$  such that  $\mathfrak{M}^2, \Gamma' \models_{pudl_2^+} \Phi$ .

**Proposition 11.**  $\mathfrak{M}^2$  has the following properties:

- (1) Both  $\leq^2$  and  $\leq_p^2$  are reflexive, transitive and connected relations.
- (2)  $w <^2 v$  iff  $w \leq^2 v$  and  $v \not\leq^2 w$
- (3)  $w <_p^2 v$  iff  $w \leq_p^2 v$  and  $v \not\leq_p^2 w$
- (4)  $R_i^2$  is an equivalence relation for each  $i \in \text{Agent}$ .
- (5) For every  $w \in W^2$ ,  $R_1^2(w) \cap \dots \cap R_n^2(w) \neq \emptyset$ .