

A Neutral Temporal Deontic STIT Logic^{*}

Kees van Berkel^(✉) and Tim Lyon

Institut für Logic and Computation, Technische Universität Wien, Austria
{kees,lyon}@logic.at

Abstract. In this work we answer a long standing request for temporal embeddings of deontic STIT logics by introducing the multi-agent STIT logic TDS. The logic is based upon atemporal utilitarian STIT logic. Yet, the logic presented here will be neutral: instead of committing ourselves to utilitarian theories, we prove the logic TDS sound and complete with respect to relational frames not employing any utilitarian function. We demonstrate how these neutral frames can be transformed into utilitarian temporal frames, while preserving validity. Last, we discuss problems that arise from employing binary utility functions in a temporal setting.

Keywords: Deontic logic · Logics of agency · Modal logic · Multi-agent STIT logic · Temporal logic · Utilitarianism

1 Introduction

With the increasing integration of automated machines in our everyday lives, the development of formal decision-making tools, which take into account moral and legal considerations, is of critical importance [2,9,10]. Unfortunately, one of the fundamental hazards of incorporating ethics into decision-making processes, is the apparent incomparability of quantitative and qualitative information—that is, moral problems most often resist quantification [16].

In contrast, utility functions are useful quantitative tools for the formal analysis of decision-making. Initially formulated in [5], the influential theory of *utilitarianism* has promoted utility calculation as a ground for *ethical deliberation*: in short, those actions generating highest utility, are the morally right actions. For this reason, utilitarianism has proven itself to be a fruitful approach in the field of formal deontic reasoning and multi-agent systems (e.g. [1,12,15]).

In particular, in the field of STIT logic—agency logics developed primarily for the formal analysis of multi-agent choice-making—the utilitarian approach has received increased attention (e.g. [1,15]). Unfortunately, each available utility function comes with its own (dis)advantages, giving rise to several puzzles (some of them addressed in [12,13]). To avoid such problems, we provide an alternative

^{*} This is a pre-print of an article published in Logic, Rationality, and Interaction. The final authenticated version is available online at: https://doi.org/10.1007/978-3-662-60292-8_25. Work funded by the projects WWTF MA16-028, FWF I2982 and FWF W1255-N23.

approach: instead of settling these philosophical issues, we develop a neutral formalism that can be appropriated to different utilitarian value assignments.

The paper's contributions can be summed up as follows: First, we provide a temporal deontic STIT logic called TDS (Sec. 2). With this logic, we answer a long standing request for temporal embeddings of deontic STIT [4,12,15]. Second, although TDS is based upon the atemporal utilitarian STIT logic from [15], the semantics of TDS will be neutral: instead of committing to utilitarianism, we prove soundness and completeness of TDS with respect to relational frames not employing any utilitarian function (Sec. 3). This approach also extends the results in [3,11,14] by showing that TDS can be characterized without using the traditional branching-time (BT+AC) structures (cf. [4]). Third, we show how neutral TDS frames can be transformed into utilitarian frames, while preserving validity (Sec. 4). Last, we discuss the philosophical ramifications of employing available utility functions in the extended, temporal setting. In particular, we will argue that binary utility assignments can turn out to be problematic.

2 A Neutral Temporal Deontic STIT Logic

In this section, we introduce the language, semantics, and axiomatization of the temporal deontic STIT logic TDS. In particular, we provide neutral relational frames characterizing the logic, which omit mention of specific utility functions. The logic will bring together atemporal deontic STIT logic, presented in [15], and the temporal STIT logic from [14].

Definition 1 (The Language \mathcal{L}_{TDS}). *Let $Ag = \{1, 2, \dots, n\}$ be a finite set of agent labels and let $Var = \{p_1, p_2, p_3, \dots\}$ be a countable set of propositional variables. The language \mathcal{L}_{TDS} is given by the following BNF grammar:*

$$\phi ::= p \mid \neg\phi \mid \phi \wedge \phi \mid \Box\phi \mid [i]\phi \mid [Ag]\phi \mid G\phi \mid H\phi \mid \otimes_i \phi$$

where $i \in Ag$ and $p \in Var$.

The logical connectives disjunction \vee , implication \rightarrow , and bi-conditional \leftrightarrow are defined in the usual way. Let \perp be defined as $p \wedge \neg p$ and define \top to be $p \vee \neg p$. The language consists of single agent STIT operators $[i]$, which are choice-operators describing that ‘agent i sees to it that’, and the grand coalition operator $[Ag]$, expressing ‘the grand coalition of agents sees to it that’. Furthermore, it contains a settledness operator \Box , which holds true of a formula that is settled true at a moment, and thus, holds true regardless of the choices made by any of the agents at that moment. The operators G and H have, respectively, the usual temporal interpretation ‘always going to be’ and ‘always has been’. Last, the operator \otimes_i expresses ‘agent i ought to see to it that’. We define \Diamond , $\langle i \rangle$, $\langle Ag \rangle$ and \ominus_i as the *duals* of \Box , $[i]$, $[Ag]$ and \otimes_i , respectively (i.e. $\Diamond\phi$ iff $\neg\Box\neg\phi$, etc.). Furthermore, let $F\phi$ iff $\neg G\neg\phi$ and $P\phi$ iff $\neg H\neg\phi$, expressing ‘ ϕ holds somewhere in the future’ and ‘ ϕ holds somewhere in the past’, respectively. Finally, deliberative STIT and deliberative ought are obtained accordingly: $[i]^d\phi$ iff $[i]\phi \wedge \Diamond\neg\phi$ and $\otimes_i^d\phi$ iff $\otimes_i\phi \wedge \Diamond\neg\phi$. For a discussion of these operators we refer to [12,14].

In line with [3,6,11,14], we provide relational frames for TDS instead of introducing the traditionally employed, BT+AC frames (cf. [4]). Explanations of the individual frame properties of Definition 2 can be found below.

Definition 2 (Relational TDS Frames and Models). *A TDS-frame is defined as a tuple $F = (W, \mathcal{R}_\square, \{\mathcal{R}_{[i]} \mid i \in Ag\}, \mathcal{R}_{[Ag]}, \mathcal{R}_G, \mathcal{R}_H, \{\mathcal{R}_{\otimes_i} \mid i \in Ag\})$. Let $\mathcal{R}_{[\alpha]}(w) := \{v \in W \mid (w, v) \in \mathcal{R}_{[\alpha]}\}$ for $[\alpha] \in \text{Boxes}$ where $\text{Boxes} := \{\square, G, H, [Ag]\} \cup \{[i] \mid i \in Ag\} \cup \{\otimes_i \mid i \in Ag\}$. Let W be a non-empty set of worlds w, v, u, \dots and:*

- For all $i \in Ag$, $\mathcal{R}_\square, \mathcal{R}_{[i]}, \mathcal{R}_{[Ag]} \subseteq W \times W$ are equivalence relations such that:
 - (C1) $\mathcal{R}_{[i]} \subseteq \mathcal{R}_\square$.
 - (C2) For all $u_1, \dots, u_n \in W$, if $\mathcal{R}_\square u_i u_j$ for all $1 \leq i, j \leq n$, then $\bigcap_i \mathcal{R}_{[i]}(u_i) \neq \emptyset$.
 - (C3) For all $w \in W$, $\mathcal{R}_{[Ag]}(w) \subseteq \bigcap_{i \in Ag} \mathcal{R}_{[i]}(w)$.
- $\mathcal{R}_G \subseteq W \times W$ is a transitive and serial binary relation and \mathcal{R}_H is the converse of \mathcal{R}_G , such that:
 - (T4) For all $w, u, v \in W$, if $\mathcal{R}_G w u$ and $\mathcal{R}_G w v$, then $\mathcal{R}_G u v$, $u = v$, or $\mathcal{R}_G v u$.
 - (T5) For all $w, u, v \in W$, if $\mathcal{R}_H w u$ and $\mathcal{R}_H w v$, then $\mathcal{R}_H u v$, $u = v$, or $\mathcal{R}_H v u$.
 - (T6) $\mathcal{R}_G \circ \mathcal{R}_\square \subseteq \mathcal{R}_{[Ag]} \circ \mathcal{R}_G$ (relation composition \circ is defined as usual).
 - (T7) For all $w, u \in W$, if $u \in \mathcal{R}_\square(w)$, then $u \notin \mathcal{R}_G(w)$.
- For all $i \in Ag$, $\mathcal{R}_{\otimes_i} \subseteq W \times W$ are binary relations such that:
 - (D8) $\mathcal{R}_{\otimes_i} \subseteq \mathcal{R}_\square$.
 - (D9) For all $w \in W$ there exists a $v \in W$ such that $\mathcal{R}_\square w v$ and for all $u \in W$, if $\mathcal{R}_{[i]} v u$ then $\mathcal{R}_{\otimes_i} w u$.
 - (D10) For all $w, v, u, z \in W$, if $\mathcal{R}_\square w v, \mathcal{R}_\square w u$ and $\mathcal{R}_{\otimes_i} u z$, then $\mathcal{R}_{\otimes_i} v z$.
 - (D11) For all $w, v \in W$, if $\mathcal{R}_{\otimes_i} w v$ then there exists $u \in W$ s.t. $\mathcal{R}_\square w u, \mathcal{R}_{[i]} u v$, and for all $z \in W$, if $\mathcal{R}_{[i]} u z$ then $\mathcal{R}_{\otimes_i} w z$.

A TDS-model is a tuple $M = (F, V)$ where F is a TDS-frame and V is a valuation mapping propositional variables to subsets of W , that is, $V: \text{Var} \rightarrow \mathcal{P}(W)$.

We label the properties of Definition 2 referring to choice **(Ci)**, those relating to temporal aspects **(Ti)**, and those capturing deontic properties **(Di)**. Observe that, since \mathcal{R}_\square is an equivalence relation, we obtain equivalence classes $\mathcal{R}_\square(w) = \{v \mid (w, v) \in \mathcal{R}_\square\}$. Furthermore, by condition **(C1)** we know that $\mathcal{R}_{[i]}$ is an equivalence relation partitioning the equivalence classes of \mathcal{R}_\square . We call $\mathcal{R}_\square(w)$ a *moment* and for each v in a moment $\mathcal{R}_\square(w)$, we refer to $\mathcal{R}_{[i]}(v)$ as a *choice-cell* for agent i at moment $\mathcal{R}_\square(w)$. In the following, we shall frequently refer to moments and choices in the above sense. Condition **(C2)** captures the pivotal *independence of agents* principle for STIT logics, ensuring that at every moment, any combination of different agents' choices is consistent: i.e., simulta-

neous choices are independent (see [4, 7C.4]). **(C3)** ensures that all agents acting together is a necessary condition for the grand coalition of agents acting.¹

The conditions on \mathcal{R}_G and \mathcal{R}_H establish that the frames we consider are irreflexive, temporal orderings of *moments*. Properties **(T4)** and **(T5)** guarantee that *histories*—i.e., maximally ordered paths of worlds passing through moments—are linear. Condition **(T6)** ensures the STIT principle of *no choice between undivided histories*: if two time-lines remain undivided at the next moment, no agent has a choice that realizes one time-line and excludes the other (see [4, 7C.3]). Consequently, this principle also ensures that the ordering of moments is linearly closed with respect to the past and allows for branching with respect to the future: in other words, TDS-frames are *treelike*.² Last, **(T7)** ensures the temporal irreflexivity of moments; i.e., the future excludes the present. For an elaborate discussion of the temporal frame conditions we refer to [14].

Last, the criteria **(D8)**–**(D11)** guarantee an essentially agentic characterization of the obligation operator \otimes_i (cf. the impartial ‘ought to be’ operator in [12]). Condition **(D8)** ensures that ideal worlds are confined to moments: i.e., the ideal worlds accessible at a moment neither lie in the future nor in the past. **(D9)** ensures that, for each agent there is at every moment a choice available that is an ideal choice (cf. the corresponding ‘ought implies can’ axiom A14). Furthermore, **(D10)** expresses that, for each agent, if a world is ideal from the perspective of a particular world at a moment, that world is ideal from the perspective of any world at that moment: i.e., ideal worlds are settled upon moments. Condition **(D11)** captures the idea that every ideal world extends to a complete ideal choice: i.e., no choice contains both ideal and non-ideal worlds. Last, note that conditions **(C2)** and **(D9)** together ensure that every combination of distinct agents’ ideal choices is consistent, i.e., non-empty.

Definition 3 (Semantics for \mathcal{L}_{TDS}). *Let M be a TDS-model and let $w \in W$ of M . The satisfaction of a formula $\phi \in \mathcal{L}_{\text{TDS}}$ in M at w is defined accordingly:*

1. $M, w \models p$ iff $w \in V(p)$
2. $M, w \models \neg\phi$ iff $M, w \not\models \phi$
3. $M, w \models \phi \wedge \psi$ iff $M, w \models \phi$ and $M, w \models \psi$
4. $M, w \models \Box\phi$ iff $\forall u \in \mathcal{R}_{\Box}(w), M, u \models \phi$
5. $M, w \models [i]\phi$ iff $\forall u \in \mathcal{R}_{[i]}(w), M, u \models \phi$
6. $M, w \models \otimes_i\phi$ iff $\forall u \in \mathcal{R}_{\otimes_i}(w), M, u \models \phi$
7. $M, w \models [Ag]\phi$ iff $\forall u \in \mathcal{R}_{[Ag]}(w), M, u \models \phi$
8. $M, w \models G\phi$ iff $\forall u \in \mathcal{R}_G(w), M, u \models \phi$
9. $M, w \models H\phi$ iff $\forall u \in \mathcal{R}_H(w), M, u \models \phi$

Global truth, validity, and semantic entailment are defined as usual (see [7]).

The axiomatization of TDS is a composition of [15], together with [14]. (Note that in the language \mathcal{L}_{TDS} each agent label represents a distinct agent.)

¹ In future work, we aim to study condition (C3) strengthened to equality, as in [14].

In such a setting, completeness is obtained by proving that each TDS-frame can be transformed into a frame (satisfying the same formulae) with strengthened (C3); hence, showing that the logic does not distinguish between the two frame classes.

² The main reason why the grand coalition operator $[Ag]$ is added to our language, is because it will allow us to axiomatize the *no choice between undivided histories* principle (see A25 of Definition 4). For a discussion of $[Ag]$ we refer to [14].

Definition 4 (Axiomatization of TDS). *For each $i \in Ag$ we have,*

<i>A0 All propositional tautologies.</i>	<i>A15 $\Diamond \otimes_i \phi \rightarrow \Box \otimes_i \phi$</i>
<i>A1 $\Box(\phi \rightarrow \psi) \rightarrow (\Box\phi \rightarrow \Box\psi)$,</i>	<i>A16 $\Box([i]\phi \rightarrow [i]\psi) \rightarrow (\otimes_i \phi \rightarrow \otimes_i \psi)$</i>
<i>A2 $\Box\phi \rightarrow \phi$</i>	<i>A17 $G(\phi \rightarrow \psi) \rightarrow (G\phi \rightarrow G\psi)$</i>
<i>A3 $\Diamond\phi \rightarrow \Box\Diamond\phi$</i>	<i>A18 $G\phi \rightarrow GG\phi$</i>
<i>A4 $[i](\phi \rightarrow \psi) \rightarrow ([i]\phi \rightarrow [i]\psi)$</i>	<i>A19 $G\phi \rightarrow F\phi$</i>
<i>A5 $[i]\phi \rightarrow \phi$</i>	<i>A20 $H(\phi \rightarrow \psi) \rightarrow (H\phi \rightarrow H\psi)$</i>
<i>A6 $\langle i \rangle \phi \rightarrow [i]\langle i \rangle \phi$</i>	<i>A21 $\phi \rightarrow GP\phi$</i>
<i>A7 $[Ag](\phi \rightarrow \psi) \rightarrow ([Ag]\phi \rightarrow [Ag]\psi)$</i>	<i>A22 $\phi \rightarrow HF\phi$</i>
<i>A8 $[Ag]\phi \rightarrow \phi$</i>	<i>A23 $FP\phi \rightarrow P\phi \vee \phi \vee F\phi$</i>
<i>A9 $\langle Ag \rangle \phi \rightarrow [Ag]\langle Ag \rangle \phi$</i>	<i>A24 $PF\phi \rightarrow P\phi \vee \phi \vee F\phi$</i>
<i>A10 $\bigwedge_{0 \leq i \leq n} \Diamond[i]\phi_k \rightarrow \Diamond \bigwedge_{0 \leq i \leq n} [i]\phi_k$</i>	<i>A25 $F\Diamond\phi \rightarrow \langle Ag \rangle F\phi$</i>
<i>A11 $\bigwedge_{1 \leq i \leq n} [i]\phi_i \rightarrow [Ag] \bigwedge_{1 \leq i \leq n} \phi_i$</i>	<i>R0 $\vdash_{TDS}(\psi \rightarrow \phi)$ and $\vdash_{TDS}\psi$ implies $\vdash_{TDS} \phi$</i>
<i>A12 $\otimes_i(\phi \rightarrow \psi) \rightarrow (\otimes_i \phi \rightarrow \otimes_i \psi)$</i>	<i>R1 $\vdash_{TDS}\phi$ implies $\vdash_{TDS}[\alpha]\phi$, $[\alpha] \in \{\Box, G, H\}$</i>
<i>A13 $\Box\phi \rightarrow ([i]\phi \wedge \otimes_i \phi)$</i>	<i>R2 $\vdash_{TDS}(\Box\neg p \wedge \Box(Gp \wedge Hp)) \rightarrow \phi$ implies</i>
<i>A14 $\otimes_i \phi \rightarrow \Diamond[i]\phi$</i>	<i>$\vdash_{TDS} \phi$, given $p \not\in \phi$</i>

A derivation of ϕ in TDS from a set Γ , written $\Gamma \vdash_{TDS} \phi$, is defined in the usual way (See [7, Def. 4.4]). When $\Gamma = \emptyset$, we say ϕ is a theorem, and write $\vdash_{TDS} \phi$.

The axioms, A1–A3, A4–A6 and A7–9 express the S5 behavior of \Box , $[i]$ (for each $i \in Ag$) and $[Ag]$, respectively. A10 is the *independence of agents* axiom. A11 captures that ‘all agents acting together implies the grand coalition of agents acting’. A13 is a bridge axiom linking \otimes_i to \Box and $[i]$ to \Box (cf. (C1) and (D8) of Definition 2). A14 corresponds to the ‘ought implies can’ principle (cf. (D9) of Definition 2). A15 ensures that, when possible, obligatory choices are settled upon moments (cf. (D10) of Definition 2). A16 can be understood as a conditional monotonicity principle for ideal choices (cf. (D11) of Definition 2). Axioms A12 and A13, together with the necessitation rule R1, ensure that \otimes_i is a normal modal operator.

With respect to the temporal axioms, A17–A19 capture the KD4 behavior of G , whereas, axioms A21 and A22 ensure that H is the converse of G . A23 and A24 capture *connectedness* of histories through moments and A25 characterizes *no choice between undivided histories*. Last, R2 is a variation of Gabbay’s irreflexivity rule (the proofs of Theorem 1 and 2 give an indication of the rule’s functions).

3 Soundness and Completeness of TDS

In this section, we prove that TDS is sound and complete relative to the class of TDS-frames. In the next section, we show how such frames are transformable into frames employing utility assignments. This allows one to model and reason about utilitarian scenarios in a more fine-grained manner, while obtaining completeness of the logic without commitment to particular utility functions.

Unless stated otherwise, all proofs in this section can be found in App. A.

Theorem 1. (SOUNDNESS OF TDS) $\forall \phi \in \mathcal{L}_{\text{TDS}}, \vdash_{\text{TDS}} \phi$ implies $\models \phi$.

We prove completeness by constructing maximal consistent sets belonging to a special class and build a canonical TDS model adopting methods from [8,14].

Definition 5. A set of formulae $\Gamma \subseteq \mathcal{L}_{\text{TDS}}$ is a maximally consistent set (MCS) iff (i) $\Gamma \not\vdash_{\text{TDS}} \perp$, and (ii) for any set $\Gamma' \subseteq \mathcal{L}_{\text{TDS}}$, if $\Gamma \subset \Gamma'$, then $\Gamma' \vdash_{\text{TDS}} \perp$.

Definition 6. (CANONICAL MODEL FOR TDS) Let $[\alpha] \in \text{Boxes}$ and let $\langle \alpha \rangle$ be the operator dual to $[\alpha]$. We define the canonical model to be the tuple $M^{dt} := (W^{dt}, \mathcal{R}_{\square}^{dt}, \{\mathcal{R}_{[i]}^{dt} \mid i \in \text{Ag}\}, \mathcal{R}_{[\text{Ag}]}^{dt}, \mathcal{R}_{\text{G}}^{dt}, \mathcal{R}_{\text{H}}^{dt}, \{\mathcal{R}_{\otimes_i}^{dt} \mid i \in \text{Ag}\}, V^{dt})$ such that:

- $W^{dt} := \{\Gamma \in \mathcal{L}_{\text{TDS}} \mid \Gamma \text{ is an MCS}\};$
- for all $\Gamma, \Delta \in W^{dt}$, $(\Gamma, \Delta) \in \mathcal{R}_{[\alpha]}^{dt}$ iff for all $\phi \in \mathcal{L}_{\text{TDS}}$, if $[\alpha]\phi \in \Gamma$, then $\phi \in \Delta$ (for each $[\alpha] \in \text{Boxes}$);
- V^{dt} is a valuation function s.t. $\forall p \in \text{Atom}, V^{dt}(p) := \{\Delta \in W^{dt} \mid p \in \Delta\}.$

Definition 7. (DIAMOND SATURATED SET [14]) Let X be a set of MCSs and let $\langle \alpha \rangle$ be dual to $[\alpha] \in \text{Boxes}$. We say that X is a diamond saturated set iff for all $\Gamma \in X$, for each $\langle \alpha \rangle \phi \in \Gamma$ there exists a $\Delta \in X$ such that $\mathcal{R}_{[\alpha]} \Gamma \Delta$ and $\phi \in \Delta$.

In order to ensure that our canonical model will be irreflexive, we introduce a mechanism that allows us to encode MCSs with information that impedes reflexive points in the model. We call these encoded sets IRR-theories and restrict our canonical model to consist of these sets only. Last, we use the notation $M|_X$ to indicate a model M whose domain is restricted to the set X (see [8, Ch.6]).

Lemma 1. Let X be a diamond saturated set with $\Gamma \in X$, $\phi \in \mathcal{L}_{\text{TDS}}$, and let $M^{dt}|_X$ be the canonical model restricted to X . Then, $M^{dt}|_X, \Gamma \models \phi$ iff $\phi \in \Gamma$.

Proof. Proven in the usual manner by induction on ϕ (see [7, Lem. 4.70]).

Following [14], we let IRR-theories be those sets of TDS formulae that (i) are maximally consistent, (ii) contain a label $\text{name}(p) := \square \neg p \wedge \square(\text{G}p \wedge \text{H}p)$, uniquely labeling a *moment* and (iii) for any world that is reachable through any ‘zig-zagging’ sequence of diamond operators, that is, every zig-zagging formula ϕ of the form,

$$\langle \alpha_1 \rangle (\phi_1 \wedge \langle \alpha_2 \rangle (\phi_2 \wedge \dots \wedge \langle \alpha_n \rangle \phi_n)) \dots$$

where $\langle \alpha_i \rangle$ is dual to $[\alpha_i] \in \text{Boxes}$ with $1 \leq i \leq n$, there exists a corresponding zig-zagging formula $\phi(q)$ (where q is a propositional variable) of the form,

$$\langle \alpha_1 \rangle (\phi_1 \wedge \langle \alpha_2 \rangle (\phi_2 \wedge \dots \wedge \langle \alpha_n \rangle (\phi_n \wedge \square \neg q \wedge \square(\text{G}q \wedge \text{H}q)))) \dots$$

labeling reachable worlds. Let us make the above formally precise:

Definition 8. (IRR-THEORY) [14] Let Zig be the set of all zig-zagging formulae in \mathcal{L}_{TDS} and let $\text{name}(p) := \square \neg p \wedge \square(\text{G}p \wedge \text{H}p)$ where p is a propositional variable. A set of formulae Γ is called an IRR-theory iff the following hold:

- Γ is a MCS and $\text{name}(p) \in \Gamma$, for some propositional variable p ;
- if $\phi \in \Gamma \cap \mathbf{Zig}$, then $\phi(q) \in \Gamma$, for some propositional variable q .

Henceforth, we refer to IRR as the set of all IRR-theories in \mathcal{L}_{TDS} .

We now present lemmata relevant to the use of IRR-theories in canonical models.

Lemma 2. *Let $\phi \in \mathcal{L}_{\text{TDS}}$ be a consistent formula. Then, there exists an IRR-theory Γ such that $\phi \in \Gamma$.*

Lemma 3. (EXISTENCE LEMMA) *Let Γ be an IRR-theory and let $\langle \alpha \rangle$ be dual to $[\alpha] \in \mathbf{Boxes}$. For each $\langle \alpha \rangle \phi \in \Gamma$ there exists an IRR-theory Δ such that $\mathcal{R}_{[\alpha]} \Gamma \Delta$.*

Subsequently, it must be shown that the canonical model *restricted* to the set IRR of IRR-theories (i.e., $M^{dt}|_{\text{IRR}}$) is in fact a TDS model (henceforth, we use W^{dt} and IRR interchangeably). First, we provide lemmata ensuring that the model satisfies the desired temporal and deontic properties of Definition 2. The first two follow from [14] and the latter four results are proven in App. A.

Lemma 4 ([14]). (PROPERTY (C2)) *Let $\Gamma_1, \dots, \Gamma_n \in \text{IRR}$ such that $\mathcal{R}_{\square}^{dt} \Gamma_i \Gamma_j$ for all $1 \leq i, j \leq n$. Then, there exists a $\Delta \in \text{IRR}$ such that $\mathcal{R}_1^{dt} \Gamma_1 \Delta, \dots, \mathcal{R}_n^{dt} \Gamma_n \Delta$.*

Lemma 5 ([14]). (PROPERTY (T6)) *Let $\Gamma, \Sigma, \Pi \in \text{IRR}$ such that $\mathcal{R}_G^{dt} \Gamma \Sigma$ and $\mathcal{R}_{\square}^{dt} \Sigma \Pi$. Then, there exists a $\Delta \in \text{IRR}$ such that $\mathcal{R}_{[Ag]}^{dt} \Gamma \Delta$ and $\mathcal{R}_G^{dt} \Delta \Pi$.*

Lemma 6. (PROPERTY (D9)) *Let $\Gamma \in \text{IRR}$. Then, there exists a $\Delta \in \text{IRR}$ such that $\mathcal{R}_{\square}^{dt} \Gamma \Delta$ and for every $\Sigma \in \text{IRR}$, if $\mathcal{R}_{[i]}^{dt} \Delta \Sigma$, then $\mathcal{R}_{\otimes_i}^{dt} \Gamma \Sigma$.*

Lemma 7. (PROPERTY (D11)) *Let $\Gamma, \Delta \in \text{IRR}$ such that $\mathcal{R}_{\otimes_i}^{dt} \Gamma \Delta$. Then, there exists a $\Sigma \in \text{IRR}$ such that $\mathcal{R}_{\square}^{dt} \Gamma \Sigma$, $\mathcal{R}_{[i]}^{dt} \Sigma \Delta$, and for all $\Pi \in \text{IRR}$, if $\mathcal{R}_{[i]}^{dt} \Sigma \Pi$, then $\mathcal{R}_{\otimes_i}^{dt} \Gamma \Pi$.*

Lemma 8. *The canonical model $M^{dt}|_{\text{IRR}}$ belongs to the class of TDS models.*

Theorem 2. (COMPLETENESS) *If $\phi \in \mathcal{L}_{\text{TDS}}$ is a consistent formula, then ϕ is satisfiable on a TDS-model.*

4 Transformations to Utilitarian Models

In this section, we investigate a truth preserving transformation from TDS models to *utilitarian* STIT models, embedded in a temporal language. In particular, we are concerned with the semantic characterization of the *dominant ought* [12, Ch.4]. We start with defining the semantic machinery needed to treat these oughts. In particular, we will introduce a utility function *util* that maps natural numbers (i.e. utilities) to worlds in our domain. In contrast to [12,15], we do not restrict the assignment of utilities to complete histories where all worlds on a maximal linear path have identical utility. The reason will be addressed at the

end of the section, where we discuss a problem related to utility assignments over histories, arising in temporal extensions of STIT.

The pivotal notion involved in the dominant ought is that of a *state*: Agent i cannot influence the choices of all other agents and, for this reason, one can regard the joint interaction of all agents excluding i , as a state (of nature) for i . To be more precise, we define a *state* $\mathcal{R}_{[i]}^s(v)$ for i at v accordingly,

$$\mathcal{R}_{[i]}^s(v) = \bigcap_{k \in Ag \setminus \{i\}} \mathcal{R}_k(v)$$

Consequently, all possible combinations of choices available to the agents $Ag \setminus \{i\}$, are the different states available at that moment to agent i .

Subsequently, we define a *preference order* \leq over choices (and subsets thereof). Let $\mathcal{R}_{[i]}(v), \mathcal{R}_{[i]}(z) \subseteq \mathcal{R}_{\square}(w)$, then weak preference is defined accordingly,

$$\mathcal{R}_{[i]}(v) \leq \mathcal{R}_{[i]}(z) \iff \forall v^* \in \mathcal{R}_{[i]}(v), \forall z^* \in \mathcal{R}_{[i]}(z), util(v^*) \leq util(z^*)$$

That is, for an agent a choice is weakly preferred over another, when all values of the possible outcomes of the former are at least as high as those of the latter (where $util(v)$ is the number assigned to v , etc). Strict preference is defined as,

$$\mathcal{R}_{[i]}(v) < \mathcal{R}_{[i]}(z) \iff \mathcal{R}_{[i]}(v) \leq \mathcal{R}_{[i]}(z) \wedge \mathcal{R}_{[i]}(z) \not\leq \mathcal{R}_{[i]}(v)$$

Next, a *dominance order* \preceq over choices $\mathcal{R}_{[i]}(v), \mathcal{R}_{[i]}(z) \subseteq \mathcal{R}_{\square}(w)$ is defined as,

$$\mathcal{R}_{[i]}(v) \preceq \mathcal{R}_{[i]}(z) \iff \forall \mathcal{R}_{[i]}^s(x) \subseteq \mathcal{R}_{\square}(w), \mathcal{R}_{[i]}(v) \cap \mathcal{R}_{[i]}^s(x) \leq \mathcal{R}_{[i]}(z) \cap \mathcal{R}_{[i]}^s(x)$$

We say an agent's choice weakly dominates another, if the values of the outcomes of the former are weakly preferred to those of the latter choice, *given any possible state available to that agent*. For a discussion of dominance orderings see [12, Ch. 4]. Again, in the usual way we obtain *strict dominance*,

$$\mathcal{R}_{[i]}(v) \prec \mathcal{R}_{[i]}(z) \iff \mathcal{R}_{[i]}(v) \preceq \mathcal{R}_{[i]}(z) \wedge \mathcal{R}_{[i]}(z) \not\preceq \mathcal{R}_{[i]}(v)$$

On the basis of the above, we now formally introduce temporal *utilitarian* STIT frames and models, defined over *relational* Kripke frames.

Definition 9 (Relational TUS Frames and Models). Let $\mathcal{R}_{[\alpha]}(w) := \{v \in W \mid (w, v) \in R_{\alpha}\}$ for $[\alpha] \in \{\square, [Ag], G, H\} \cup \{[i] \mid i \in Ag\}$. A *relational Temporal Utilitarian STIT frame (TUS-frame)* is defined as a tuple $F = (W, \mathcal{R}_{\square}, \{\mathcal{R}_{[i]} \mid i \in Ag\}, \mathcal{R}_{[Ag]}, \mathcal{R}_G, \mathcal{R}_H, util)$ where W is a non-empty set of worlds w, v, u, \dots and:

- For all $i \in Ag$, $\mathcal{R}_{\square}, \mathcal{R}_{[i]}, \mathcal{R}_{[Ag]} \subseteq W \times W$ are equivalence relations for which conditions **(C1)**–**(C3)** of Definition 2 hold.
- $\mathcal{R}_G \subseteq W \times W$ is a transitive and serial binary relation, whereas \mathcal{R}_H is the converse of \mathcal{R}_G , and the conditions **(T4)**–**(T7)** of Definition 2 hold.
- $util : W \mapsto \mathbb{N}$ is a utility function assigning each world in W to a natural.

A TUS-model is a tuple $M = (F, V)$ where F is a TUS-frame and V is a valuation function assigning propositional variables to subsets of W : i.e., $V : Var \mapsto \mathcal{P}(W)$.

Notice that the above TUS frames only differ from TDS frames through replacing the relations \mathcal{R}_{\otimes_i} and corresponding conditions (D8)-(D11) (for each $i \in Ag$) with the utility function $util$. We observe that the assignment of utilities to worlds is agent-independent. Nevertheless, since the choices of an agent depend on which worlds are inside the choice-cells available to the agent, the resulting obligations are in fact agent-dependent. Let us define the new semantics:

Definition 10 (Semantics of TUS models). *Let M be a TUS-model, $w \in W$ of M and $\|\phi\|_M = \{w \mid M, w \models \phi\}$. We define satisfaction of a formula $\phi \in \mathcal{L}_{\text{TDS}}$ as follows:*

- Clause (1)-(10) are the same as those from Definition 3, with the exception of clause (7), which we replace by the following clause (7*):

$$M, w \models \otimes_i \phi \text{ iff } \forall \mathcal{R}_{[i]}(v) \subseteq \mathcal{R}_{\square}(w) \text{ if } \mathcal{R}_{[i]}(v) \not\subseteq \|\phi\| \text{ then } \exists \mathcal{R}_{[i]}(z) \subseteq \mathcal{R}_{\square}(w) \text{ s.t.} \\ (i) \mathcal{R}_{[i]}(v) \prec \mathcal{R}_{[i]}(z), (ii) \mathcal{R}_{[i]}(z) \subseteq \|\phi\| \text{ and} \\ (iii) \forall \mathcal{R}_{[i]}(x) \subseteq \mathcal{R}_{\square}(w), \mathcal{R}_{[i]}(z) \preceq \mathcal{R}_{[i]}(x) \text{ implies } \mathcal{R}_{[i]}(x) \subseteq \|\phi\|$$

Clause (7*) is interpreted accordingly: Agent i ought to see to it that ϕ iff for every choice $\mathcal{R}_{[i]}(v)$ available to i that does not guarantee ϕ there (i) exists a strictly dominating choice $\mathcal{R}_{[i]}(z)$ that (ii) does guarantee ϕ and (iii) every weakly dominating choice $\mathcal{R}_{[i]}(x)$ over $\mathcal{R}_{[i]}(z)$ also guarantees ϕ . In other words, all choices not guaranteeing ϕ are strictly dominated only by choices guaranteeing ϕ . (We note that clause (7*) is obtained through an adaption of the definition provided in [12] to relational frames.) We show that the logic TDS is also sound and complete with respect to the class of TUS-frames.

Theorem 3. (SOUNDNESS) $\forall \phi \in \mathcal{L}_{\text{TDS}}, \text{ if } \vdash_{\text{TDS}} \phi, \text{ then } \mathcal{C}_f^u \models \phi.$

Proof. We prove by induction on the given derivation of ϕ in TDS. The argument for axioms A0-A6 and A12 is the same as in Theorem 1. The validity of the axioms A7-A11 can be easily checked by applying semantic clause (7*) of Definition 9.

We now prove that the class \mathcal{C}_f^u of TUS-frames characterizes the same set of formulae as the class \mathcal{C}_f^d of TDS frames. We prove both directions separately:

Theorem 4. $\forall \phi \in \mathcal{L}_{\text{TDS}} \text{ we have } \mathcal{C}_f^u \models \phi \text{ implies } \mathcal{C}_f^d \models \phi.$

Proof. We prove by contraposition assuming $\mathcal{C}_f^d \not\models \phi$. Hence, there is a TDS-model, $\mathcal{M}^d = (\mathcal{W}, \mathcal{R}_{\square}, \{\mathcal{R}_i \mid i \in Ag\}, \mathcal{R}_H, \mathcal{R}_G, \mathcal{R}_{Ag}, \{\mathcal{R}_{\otimes_i} \mid i \in Ag\}, \mathcal{V})$ such that $\mathcal{M}^d, w \models \neg\phi$ for some $w \in \mathcal{W}$. We use \mathcal{M}^d to construct a model M in \mathcal{C}_f^u , such that:

$$M = (\mathcal{W}, \mathcal{R}_{\square}, \{\mathcal{R}_i \mid i \in Ag\}, \mathcal{R}_G, \mathcal{R}_H, \mathcal{R}_{Ag}, util, \mathcal{V})$$

We show that $M, w' \models \neg\phi$ for some $w' \in \mathcal{W}$. To define M let $\mathcal{W} := \mathcal{W}$, $\mathcal{R}_{\square} := \mathcal{R}_{\square}$, $\mathcal{R}_i := \mathcal{R}_i$, $\mathcal{R}_H := \mathcal{R}_H$, $\mathcal{R}_G := \mathcal{R}_G$, $\mathcal{R}_{Ag} := \mathcal{R}_{Ag}$, $\mathcal{V}(p) := \mathcal{V}(p)$ and let $util$ be a function assigning each $w \in \mathcal{W}$ to a natural number, satisfying the following criteria:

1. $\forall i \in Ag, \forall w, v, z \in \mathcal{W}$, if $v, z \in \mathcal{R}_\square(w)$, $v \in \mathcal{R}_i^s(w) \setminus \mathcal{R}_{\otimes_i}(w)$, and $z \in \mathcal{R}_i^s(w) \cap \mathcal{R}_{\otimes_i}(w)$, then $\text{util}(v) \leq \text{util}(z)$;
2. $\forall w, v, z \in \mathcal{W}$, if $v \in \mathcal{R}_\square(w) \setminus \mathcal{R}_{\otimes_{Ag}}(w)$ and $z \in \mathcal{R}_{\otimes_{Ag}}(w)$, then $\text{util}(v) < \text{util}(z)$;
3. $\forall w, u, z \in \mathcal{W}$, if $v, z \in \mathcal{R}_i^s(w) \cap \mathcal{R}_{\otimes_i}(w)$, then $\text{util}(v) = \text{util}(z)$;

Let $\mathcal{R}_{\otimes_{Ag}} := \bigcap_{i \in Ag} \mathcal{R}_{\otimes_i}$, we call $\mathcal{R}_{[i]}(v) \subseteq \mathcal{R}_{\otimes_i}(w)$ an *optimal choice* for agent i . (It can be easily checked that the function util can be constructed.)

We state the following useful lemma (the proof of which is found in App. A).

Lemma 9. *The following holds for any TDS frame:*

- (1) $\forall v \in \mathcal{R}_\square(w), \mathcal{R}_\square(w) = \mathcal{R}_\square(v)$; (2) $\forall v \in \mathcal{R}_i(w), \mathcal{R}_i(w) = \mathcal{R}_i(v)$;
- (3) $\forall v \in \mathcal{R}_i^s(w), \mathcal{R}_i^s(w) = \mathcal{R}_i^s(v)$; (4) $\forall v \in \mathcal{R}_\square(w)$ we get $\mathcal{R}_{\otimes_i}(v) = \mathcal{R}_{\otimes_i}(w)$;
- (5) $\forall \mathcal{R}_{[i]}(z) \subseteq \mathcal{R}_\square(w)$, either $\mathcal{R}_{[i]}(z) \subseteq \mathcal{R}_{\otimes_i}(w)$ or $\mathcal{R}_{[i]}(z) \cap \mathcal{R}_{\otimes_i}(w) = \emptyset$.

We observe that conditions **(C1)**–**(C3)** and **(T4)**–**(T7)** will be satisfied in \mathbf{M} since all of the relations of \mathcal{M}^d , with the exception of \mathcal{R}_{\otimes_i} , are identical to those in \mathbf{M} . Moreover, util complies with Definition 9 and so \mathbf{M} is in fact a TUS model. The desired claim will follow if we additionally show that $\forall \psi \in \mathcal{L}_{\text{TDS}}$ and $\forall w \in \mathcal{W}$:

$$\mathcal{M}^d, w \models \psi \iff \mathbf{M}, w \models \psi$$

We prove the claim by induction on the complexity of ψ .

Base Case. Let ψ be a propositional variable p . By the definition of \mathbf{V} in \mathbf{M} it follows directly that $\mathcal{M}^d, w \models p$ iff $w \in \mathcal{V}$ iff $w \in \mathbf{V}$ iff $\mathbf{M}, w \models p$.

Inductive Step. The cases for the propositional connectives and the modalities $[\alpha] \in \{\square, \mathbf{H}, \mathbf{G}, [Ag]\} \cup \{[i] \mid i \in Ag\}$ are straightforward. We consider the non-trivial case when ψ is of the form $\otimes_i \phi$. Let us first prove the left to right direction.

(\implies) Assume $\mathcal{M}^d, w \models \otimes_i \phi$. We show that $\mathbf{M}, w \models \otimes_i \phi$. By the semantics for \otimes_i (Definition 9) it suffices to prove that: $\forall \mathcal{R}_i(v) \subseteq \mathcal{R}_\square(w)$ if $\mathcal{R}_i(v) \not\subseteq \|\phi\|_{\mathbf{M}}$, then $\exists \mathcal{R}_i(u) \subseteq \mathcal{R}_\square(w)$ such that the following three clauses hold: (i) $\mathcal{R}_i(v) \prec \mathcal{R}_i(u)$; (ii) $\mathcal{R}_i(u) \subseteq \|\phi\|_{\mathbf{M}}$; and (iii) $\forall \mathcal{R}_i(x) \subseteq \mathcal{R}_\square(w)$, $\mathcal{R}_i(u) \preceq \mathcal{R}_i(x)$ implies $\mathcal{R}_i(x) \subseteq \|\phi\|_{\mathbf{M}}$.

Let $\mathcal{R}_i(v) \subseteq \mathcal{R}_\square(w)$ be arbitrary and assume that $\mathcal{R}_i(v) \not\subseteq \|\phi\|_{\mathbf{M}}$. We prove that there is a $\mathcal{R}_i(u) \subseteq \mathcal{R}_\square(w)$ for which conditions (i)–(iii) hold. First, we prove the existence of such a $\mathcal{R}_i(u) \subseteq \mathcal{R}_\square(w)$: By **(C1)** and **(D9)** of Definition 2, we know,

$$\exists u \in \mathcal{W} \text{ such that } \mathcal{R}_i(u) \subseteq \mathcal{R}_\square(w) \text{ and } \mathcal{R}_i(u) \subseteq \mathcal{R}_{\otimes_i}(w). \quad (1)$$

We also know by **(D9)** that $\forall j \in Ag \setminus \{i\}, \exists u_j \in \mathcal{R}_\square(w)$ such that $\mathcal{R}_j(u_j) \subseteq \mathcal{R}_{\otimes_j}(w)$. By **(IOA)** we know that $\bigcap_{j \in Ag \setminus \{i\}} \mathcal{R}_j(u_j) \cap \mathcal{R}_i(u) \neq \emptyset$, i.e., there exists a $u^* \in \bigcap_{j \in Ag \setminus \{i\}} \mathcal{R}_j(u_j) \cap \mathcal{R}_i(u)$. Consequently, we obtain the following statement,

$$u^* \in \bigcap_{j \in Ag \setminus \{i\}} \mathcal{R}_{\otimes_j}(w) \cap \mathcal{R}_{\otimes_i}(w) = \mathcal{R}_{\otimes_{Ag}}(w). \quad (2)$$

Last, by construction of \mathbf{M} we know $\mathcal{R}_i(u) = \mathcal{R}_i(u)$. We show that (i)–(iii) hold:

(i) We show $\mathcal{R}_i(v) \prec \mathcal{R}_i(u)$, that is, (a) $\mathcal{R}_i(v) \preceq \mathcal{R}_i(u)$ and (b) $\mathcal{R}_i(u) \not\preceq \mathcal{R}_i(v)$:

(a) Recall, $\mathcal{R}_i(v) \not\subseteq \|\phi\|_{\mathbf{M}}$, we know $\exists v^* \in \mathcal{R}_i(v)$ s.t. $\mathbf{M}, v^* \not\models \phi$. By definition of \mathbf{M} , $v^* \in \mathcal{R}_i(v)$ and by (IH) we get $\mathcal{M}^d, v^* \not\models \phi$. Consequently, by the assumption

that $\mathcal{M}^d, w \models \otimes_i \phi$, and the fact that $\mathcal{M}^d, v^* \not\models \phi$, it follows that $v^* \notin \mathcal{R}_{\otimes_i}(w)$. Hence, we know that $\mathcal{R}_i(v) \not\subseteq \mathcal{R}_{\otimes_i}(w)$, which implies $\mathcal{R}_{\otimes_i}(w) \cap \mathcal{R}_i(v) = \emptyset$ by Lemma 9–(5). Therefore, by this fact along with statement (1) above, we know that,

For all $x, u^\nabla, v^\nabla \in \mathcal{W}$, if $v^\nabla \in \mathcal{R}_i^s(x) \cap \mathcal{R}_i(v)$ and $u^\nabla \in \mathcal{R}_i^s(x) \cap \mathcal{R}_i(u)$, then $v^\nabla \in \mathcal{R}_i^s(x) \setminus \mathcal{R}_{\otimes_i}(w)$ and $u^\nabla \in \mathcal{R}_i^s(x) \cap \mathcal{R}_{\otimes_i}(w)$.

Let $x, u^\nabla, v^\nabla \in \mathcal{W}$ be arbitrary and assume that $v^\nabla \in \mathcal{R}_i^s(x) \cap \mathcal{R}_i(v)$ and $u^\nabla \in \mathcal{R}_i^s(x) \cap \mathcal{R}_i(u)$. By the statement above, it follows that $v^\nabla \in \mathcal{R}_i^s(x) \setminus \mathcal{R}_{\otimes_i}(w)$ and $u^\nabla \in \mathcal{R}_i^s(x) \cap \mathcal{R}_{\otimes_i}(w)$, which in conjunction with criterion 1 on the function util implies that $\text{util}(v^\nabla) \leq \text{util}(u^\nabla)$. Therefore, the following holds,

For all $x, u^\nabla, v^\nabla \in \mathcal{W}$, if $v^\nabla \in \mathcal{R}_i^s(x) \cap \mathcal{R}_i(v)$ and $u^\nabla \in \mathcal{R}_i^s(x) \cap \mathcal{R}_i(u^\nabla)$, then $\text{util}(v^\nabla) \leq \text{util}(u)$.

It follows that $\forall \mathcal{R}_i^s(x) \subseteq \mathcal{R}_\square(w)$, $\mathcal{R}_i^s(x) \cap \mathcal{R}_i(v) \leq \mathcal{R}_i^s(x) \cap \mathcal{R}_i(u)$. Hence, by the definition of \preceq and the definition of \mathbf{M} , we obtain $\mathbf{R}_i(v) \preceq \mathbf{R}_i(u)$.

(b) We need to show $\mathbf{R}_i(u) \not\preceq \mathbf{R}_i(v)$. By definition of \preceq , it suffices to show that $\exists x, \exists u^\nabla, \exists v^\nabla \in \mathcal{W}$ s.t. $\mathbf{R}_i(x) \subseteq \mathbf{R}_\square(w)$, $u^\nabla \in \mathbf{R}_i(u) \cap \mathbf{R}_i^s(x)$, $v^\nabla \in \mathbf{R}_i(v) \cap \mathbf{R}_i^s(x)$ and $\text{util}(v^\nabla) < \text{util}(u^\nabla)$. Consider $\bigcap_{j \in \text{Ag} \setminus i} \mathcal{R}_j(u_j) \cap \mathcal{R}_i(u) \neq \emptyset$ from statement (2). Let $\mathbf{R}_i^s(x) := \bigcap_{j \in \text{Ag} \setminus i} \mathcal{R}_j(u_j)$. Clearly, $\mathbf{R}_i^s(x) \subseteq \mathbf{R}_\square(w)$. By (IOA) we know that $\mathcal{R}_i^s(x) \cap \mathcal{R}_i(v) \neq \emptyset$ (where $\mathcal{R}_i^s(x) = \bigcap_{j \in \text{Ag} \setminus i} \mathcal{R}_j(u_j)$), and so, $\mathbf{R}_i^s(x) \cap \mathbf{R}_i(v) \neq \emptyset$ by the definition of \mathbf{M} . Therefore, $\exists v^\nabla \in \mathbf{R}_i^s(x) \cap \mathbf{R}_i(v)$. Since $u^* \in \bigcap_{j \in \text{Ag} \setminus i} \mathcal{R}_j(u_j) \cap \mathcal{R}_i(u)$ (see paragraph above statement (2)), we know that $u^* \in \bigcap_{j \in \text{Ag} \setminus i} \mathcal{R}_j(u_j) \cap \mathbf{R}_i(u)$, implying that $u^* \in \mathbf{R}_i^s(x) \cap \mathbf{R}_i(u)$. Since also $\mathcal{R}_i(v) \cap \mathcal{R}_{\otimes_{\text{Ag}}}(w) = \emptyset$, as derived in part (i), we obtain $v^\nabla \in \mathbf{R}_\square(w) \setminus \mathcal{R}_{\otimes_{\text{Ag}}}(w)$. By criterion 2 of util , and the facts $v^\nabla \in \mathbf{R}_\square(w) \setminus \mathcal{R}_{\otimes_{\text{Ag}}}(w)$ and $u^* \in \mathcal{R}_{\otimes_{\text{Ag}}}(w)$, by statement (2), we have that $\text{util}(v^\nabla) < \text{util}(u^*)$. Therefore, $\mathbf{R}_i(u) \not\preceq \mathbf{R}_i(v)$.

(ii) By assumption $\mathcal{R}_{\otimes_i}(w) \subseteq \|\phi\|_{\mathcal{M}^d}$ and statement (1) we get $\mathcal{R}_i(u) \subseteq \mathcal{R}_{\otimes_i}(w)$. By IH we have $\|\phi\|_{\mathcal{M}^d} = \|\phi\|_{\mathbf{M}}$ and since $\mathcal{R}_i(u) = \mathbf{R}_i(u)$ we know $\mathbf{R}_i(u) \subseteq \|\phi\|_{\mathbf{M}}$.

(iii) We prove the case by contraposition and show that $\forall \mathbf{R}_i(x) \subseteq \mathbf{R}_\square(w)$, if $\mathbf{R}_i(x) \not\subseteq \|\phi\|$, then $\mathbf{R}_i(u) \not\preceq \mathbf{R}_i(x)$. Let $\mathbf{R}_i(x)$ be an arbitrary choice-cell in $\mathbf{R}_\square(w)$ and assume that $\mathbf{R}_i(x) \not\subseteq \|\phi\|_{\mathbf{M}}$. We aim to prove that $\mathbf{R}_i(u) \not\preceq \mathbf{R}_i(x)$. By definition of \preceq it suffices to show that $\exists \mathbf{R}_i^s(y) \subseteq \mathbf{R}_\square(w)$ such that $\exists u^\nabla \in \mathbf{R}_i(u) \cap \mathbf{R}_i^s(y)$, $\exists x^\nabla \in \mathbf{R}_i(x) \cap \mathbf{R}_i^s(y)$, and $\text{util}(x^\nabla) < \text{util}(u^\nabla)$.

By the assumption that $\mathbf{R}_i(x) \not\subseteq \|\phi\|_{\mathbf{M}}$, we know $\exists x^\nabla \in \mathbf{R}_i(x)$ such that $\mathbf{M}, x^\nabla \not\models \phi$. Clearly, $x^\nabla \in \mathcal{R}_i(x)$, and by (IH) we know that $\mathcal{M}^d, x^\nabla \not\models \phi$. Since $\mathcal{M}^d, w \models \otimes_i \phi$, we obtain $(w, x^\nabla) \notin \mathcal{R}_{\otimes_i}$, and by Lemma 9–(5) we obtain $\mathcal{R}_i(x) \not\subseteq \mathcal{R}_{\otimes_i}(w)$.

By statement (2) we had $u^* \in \mathcal{R}_{\otimes_{\text{Ag}}}(w)$ and $u^* \in \mathcal{R}_{\otimes_i}(w)$. Also, we know $u^* \in \mathcal{R}_i(u)$ by paragraph preceding statement (2). Since, $u^* \in \bigcap_{j \in \text{Ag} \setminus \{i\}} \mathcal{R}_j(u_j) \cap \mathcal{R}_i(u)$, we also have $u^* \in \bigcap_{j \in \text{Ag} \setminus \{i\}} \mathcal{R}_j(u_j)$. Let $\mathcal{R}_i^s(u^*) := \bigcap_{j \in \text{Ag} \setminus \{i\}} \mathcal{R}_j(u_j)$. By (IOA) we obtain $\mathcal{R}_i(x) \cap \mathcal{R}_i^s(u^*) \neq \emptyset$, implying that there exists some $x^\nabla \in \mathcal{R}_i(x) \cap \mathcal{R}_i^s(u^*)$. It follows from (D9) and the fact $\mathcal{R}_i(x) \not\subseteq \mathcal{R}_{\otimes_i}(w)$ that $x^\nabla \notin \mathcal{R}_{\otimes_{\text{Ag}}}(w)$, which with the fact $u^* \in \mathcal{R}_{\otimes_{\text{Ag}}}(w)$, implies by definition of util (criterion 2) that $\text{util}(x^\nabla) < \text{util}(u^*)$. By the definition of \mathbf{M} , we have

$x^\nabla \in R_i(x) \cap R_i^s(u^*)$, $u^* \in R_i(u) \cap R_i^s(u^*)$ and $\text{util}(x^\nabla) < \text{util}(u^*)$, which implies the desired claim.

(\Leftarrow) We now prove the right to left direction: Assume $M, w \models \otimes_i \phi$. We reason towards a contradiction by assuming $M^d, w \not\models \otimes_i \phi$. Hence, there exists a world $v \in \mathcal{R}_{\otimes_i}(w)$ such that $M^d, v \not\models \phi$. By (D11) we obtain $\mathcal{R}_{[i]}(v) \subseteq \mathcal{R}_{\otimes_i}(w)$ and hence $\mathcal{R}_{[i]}(v) \not\subseteq \|\phi\|_{M^d}$. By (IH) and the definition of M , we obtain $R_i(v) \not\subseteq \|\phi\|_M$. This fact, in conjunction with the assumption $M, w \models \otimes_i \phi$, implies that there exists some $R_i(z) \subseteq R_\square(w)$ such that the following holds: (i) $R_i(v) \prec R_i(z)$; (ii) $R_i(z) \subseteq \|\phi\|_M$; and (iii) $\forall R_i(x) \subseteq R_\square(w)$, $R_i(z) \preceq R_i(x)$ implies $R_i(x) \subseteq \|\phi\|_M$.

By Lemma 9–(5) and the fact that $R_i(z) = \mathcal{R}_i(z)$, we know that either (a) $\mathcal{R}_i(z) \subseteq \mathcal{R}_{\otimes_i}(w)$ holds or (b) $\mathcal{R}_i(z) \cap \mathcal{R}_{\otimes_i}(w) = \emptyset$ holds.

Assume (a). We know $R_i(v) \prec R_i(z)$ and therefore, $R_i(z) \not\preceq R_i(v)$. Hence, $\exists R_i^s(x) \subseteq R_\square(w)$, $\exists z^* \in R_i(z) \cap R_i^s(x)$, $\exists v^* \in R_i(v) \cap R_i^s(x)$ such that $\text{util}(v^*) < \text{util}(z^*)$. We also know $\mathcal{R}_i(v) \subseteq \mathcal{R}_{\otimes_i}(w)$ and $\mathcal{R}_i(z) \subseteq \mathcal{R}_{\otimes_i}(w)$ and thus we obtain $z^*, v^* \in \mathcal{R}_{\otimes_i} \cap \mathcal{R}_i^s(x)$. Consequently, by the definition of util (criterion 3), we get $\text{util}(v^*) = \text{util}(z^*)$. Contradiction.

Assume (b). We know $R_i(v) \prec R_i(z)$ and therefore, $R_i(z) \not\preceq R_i(v)$. Hence, $\exists R_i^s(x) \subseteq R_\square(w)$, $\exists z^* \in R_i(z) \cap R_i^s(x)$, $\exists v^* \in R_i(v) \cap R_i^s(x)$ such that $\text{util}(z^*) \not\leq \text{util}(v^*)$. Then, by definition of util (criterion 1), either (I) $z^* \notin \mathcal{R}_i^s(x) \setminus \mathcal{R}_{\otimes_i}(w)$ or (II) $v^* \notin \mathcal{R}_i^s(x) \cap \mathcal{R}_{\otimes_i}(w)$. Suppose (I), since $z^* \in R_i^s(x)$ we infer $z^* \in \mathcal{R}_i^s(x)$ and thus conclude $z^* \in \mathcal{R}_{\otimes_i}(w)$. However, by earlier assumption $\mathcal{R}_i(z) \cap \mathcal{R}_{\otimes_i}(w) = \emptyset$ we obtain $z^* \notin \mathcal{R}_{\otimes_i}(w)$. Contradiction. Suppose (II), then since $v^* \in \mathcal{R}_i^s(x)$ we infer $v^* \notin \mathcal{R}_{\otimes_i}(w)$. However, $\mathcal{R}_{[i]}(v) \subseteq \mathcal{R}_{\otimes_i}(w)$. Contradiction.

Corollary 1. (COMPLETENESS) $\forall \phi \in \mathcal{L}_{\text{TDS}}$, if $\mathcal{C}_f^u \models \phi$, then $\vdash_{\text{TDS}} \phi$.

Proof. Follows from Theorem 4 above, together with Theorem 2.

Theorem 5. $\forall \phi \in \mathcal{L}_{\text{TDS}}$, we get $\mathcal{C}_f^d \models \phi$ implies $\mathcal{C}_f^u \models \phi$.

Proof. Follows from Theorem 2 together with Theorem 3.

The Problem with Two-Valued Utility Functions. A well studied candidate function for assigning utilities to *histories*, is the *two-valued* approach where the range of utilities is $\{0, 1\}$ (e.g. [12, 15]). As a concluding remark of the present section, we briefly discuss the philosophical ramifications of using binary utility functions in a temporal setting.

Observe that, at a moment where all worlds have a utility of 1 (or all 0), every obligation becomes vacuously satisfied by definition—in such a scenario we would have $\otimes_i \phi$ iff $\square \phi$ —and every choice for each agent will ensure all optimal outcomes (see clause (7*) of Definition 10).³ If in such a scenario, following [12, 15], utilities are assigned to complete histories and thus remain constant through time, all obligations will also be vacuously satisfied at every future moment from thereon (namely, as one moves into the future, the set of histories passing through a moment can only decrease or stay the same). That at such

³ This also holds when all intersections of choices of agents contain both a 1 and a 0.

moments all obligations are vacuously satisfied means that no obligation can be violated. Unfortunately, this also implies that at such moments *contrary-to-duty* (CTD) reasoning—i.e., reasoning about obligations that come into being when a previous obligation has been violated—becomes impossible because CTD obligations require the possibility to violate one’s obligations in the first place (e.g. see [17]).

In order to reason with CTD obligations in *temporal* utilitarian STIT logics, we need to ensure that obligations can be violated, that is, we must consider deliberative obligations: $\otimes_i^d \phi := \otimes_i \phi \wedge \neg \Box \phi$. This means that, for an obligation $\otimes_i^d \phi$ to hold, there exists a choice that does not guarantee ϕ and, by definition, the latter choice must be strictly dominated by (only) ϕ choices. In the binary setting this means that for all optimal choices, there is at least one outcome with a strictly higher utility (which must be 1). Unfortunately, this has a drawback since at such moments *at least* one of the following holds: (1) Worlds in the intersection of all agents acting in accordance with their duty *all have value* 1. (2) Worlds in the intersection of all agents violating their duty *all have value* 0.

Relative to the aforementioned, Fig. 1 illustrates the (only) three scenarios possible in a two-agents, two-choices setting: Sub-figure (i) implies the impossibility of future CTD reasoning in all cases in which at least one agent satisfies its obligation. Sub-figure (ii) implies that there is no future CTD possible in every case witnessing at least one agent violating its obligation. Last, sub-figure (iii) indicates that future CTD obligations can only occur if one of the agents satisfies her obligation if and only if the other violates his. (With the impossibility of future CTD reasoning we mean that from that moment onward, all obligations will be vacuously satisfied.) All three cases are undesirable since they do not allow for future recuperation in those situations in which they clearly should.

The above exhibits that, although \otimes_i does not depend on any temporal aspect (e.g. [15]), we can identify utility functions that are less suitable for temporal extensions of STIT. Binary functions relative to moments only, do not cause these problems, although they have their own issues [12]. In the case where the function ranges over the set of reals, it is possible to assign utilities in such a way that there is always CTD reasoning possible. In future work, we aim to specify such utility functions, making particular use of temporal aspects of TDS-frames.

5 Conclusion and Future Work

In this paper, we extended deontic STIT logic [15] to the temporal setting, incorporating the logic from [14]. In doing so, we answered a long standing open question for temporal embeddings of deontic STIT (e.g. see [4,12,15]). We showed that the resulting logic TDS is sound and complete with respect to its class of frames. We dubbed these frames *neutral* since they allowed us to obtain adequacy of the calculus, while allowing us to refrain from committing to specific utility functions. Subsequently, we showed how these neutral frames can be transformed into particular utilitarian models, while preserving truth. We also briefly argued

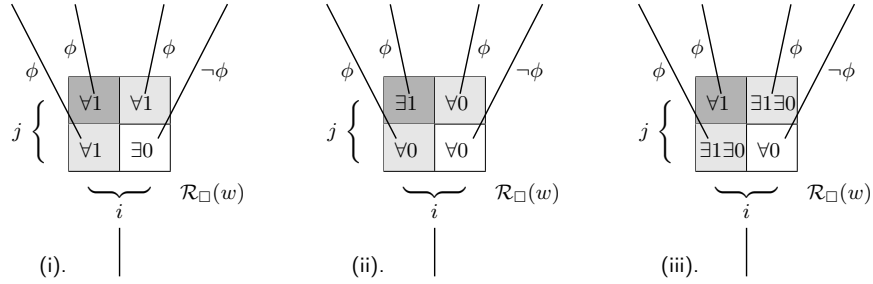


Fig. 1. The only three scenarios where $\otimes_i \phi \wedge \otimes_j \phi \wedge \neg \square \phi$ holds true at $\mathcal{R}_\square(w)$ (for $Ag = \{i, j\}$ with 2 choices). Choices of i are vertically presented, those of j horizontally. The symbol $\forall n$ means every history is assigned value n , and $\exists n$ means that some history is assigned n , for $n \in \{0, 1\}$. Optimal choices are shaded and darker shaded when overlapping. At all $\forall k$ outcomes (with $k \in \{0, 1\}$), CTD reasoning becomes impossible.

that in a temporal setting, binary value assignments to histories can generate undesirable behavior with respect to contrary-to-duty obligations.

For future work, we leave open the problem of whether temporal STIT (from [14]) and its deontic extension TDS are decidable. Furthermore, we aim to investigate alternative utility assignments that explicitly exploit the temporal aspects of TDS; e.g., it might be interesting to consider a dynamic approach taking into account that natural agents have limited foresight relative to (future) utilities.

References

1. Abarca, A.I.R., Broersen, J.: A Logic of Objective and Subjective Oughts. In: JELIA 2019: Joint European Conference on Logics in Artificial Intelligence. Springer, Cham, pp.629-641 (2019)
2. Arkoudas, K., Bringsjord S., Bello, P.: Toward ethical robots via mechanized deontic logic. In: AAAI Fall Symposium on Machine Ethics, pp.17-23 (2005)
3. Balbiani, P., Herzig, A., Troquard, N.: Alternative axiomatics and complexity of deliberative STIT theories. *Journal of Philosophical Logic*, 37(4), pp.387-406. Springer (2008)
4. Belnap, N, Perloff, M., Xu, M.: Facing the future: Agents and choices in our indeterminist world. Oxford University Press on Demand, Oxford (2001)
5. Bentham, J.: An Introduction to the Principles of Morals and Legislation. (1789)
6. Berkel, K. van, Lyon, T.: Cut-free Calculi and Relational Semantics for Temporal STIT Logics. In: JELIA 2019: Joint European Conference on Logics in Artificial Intelligence, Springer Cham (2019)
7. Blackburn, P., de Rijke, M., Venema, Y.: Modal logic. Cambridge University Press, Cambridge (2001)
8. Gabbay, D. M., Hodkinson, I., Reynolds, M.: Temporal logic: Mathematical foundations and computational aspects. Oxford University Press, Oxford (1994)
9. Gerdes, J.C., Thornton, S.M.: Implementable ethics for autonomous vehicles. In: *Autonomes fahren*, pp.87-102. Springer Vieweg, Berlin, Heidelberg (2015)

10. Goodall, N.J.: Machine ethics and automated vehicles. In: Road vehicle automation, pp.93–102 . Springer, Cham (2014)
11. Herzig, A., Schwarzenrüber, F.: Properties of logics of individual and group agency. In: Advances in Modal Logic (7), pp. 133–149. College Publications (2008)
12. Horty, J.: Agency and Deontic Logic. Oxford University Press (2001)
13. Horty, J. and Pacuit, E.: Action Types in STIT Semantics. The Review of Symbolic Logic 10(4), pp. 617–637 (2017)
14. Lorini, E.: Temporal STIT logic and its application to normative reasoning. Journal of Applied Non-Classical Logics 23 (4), pp. 372–399 (2013)
15. Murakami, Y.: Utilitarian deontic logic. In: Advances in Modal Logic (5), pp. 211–230. King’s College Publications (2005)
16. Nayeypour, M. and Koehn, D.: The Ethics of Quality: Problems and Preconditions. In: Journal of Business Ethics 44(1), pp. 37–48. Kluwer Academic Publishers (2003)
17. Prakken, H., Sergot, M.: Contrary-to-duty obligations. Studia Logica 57(1), pp.91–115 (1996)

A Proofs

Theorem 1 (SOUNDNESS) $\forall \phi \in \mathcal{L}_{\text{TDS}}, \vdash_{\text{TDS}} \phi$ implies $\models \phi$.

Proof. It suffices to show that all axioms are valid and all inference rules preserve validity over the class of TDS frames. The rules $R0$ and $R1$, as well as axioms $A0$ – $A11$, and $A17$ – $A25$ can be easily checked (See [14]). We show that $A13$ – $A16$ are valid and that the $R2$ preserves validity. Let M be an arbitrary TDS-model with w a world in M .

A13. Assume $M, w \models \Box \phi$ and also that $\mathcal{R}_{[i]}wu$ and $\mathcal{R}_{\otimes_i}wv$. By conditions (C1) and (D8), we know that $\mathcal{R}_{[i]} \subseteq \mathcal{R}_{\Box}$ and $\mathcal{R}_{\otimes_i} \subseteq \mathcal{R}_{\Box}$, respectively. Therefore, it follows that $\mathcal{R}_{\Box}wu$ and $\mathcal{R}_{\Box}wv$, which implies $M, u \models \phi$ and $M, v \models \phi$ by the assumption. This implies that $M, w \models [i]\phi$ and $M, w \models \otimes_i \phi$.

A14. Assume $M, w \models \otimes_i \phi$. By condition (D9), there exists a v such that $\mathcal{R}_{\Box}wv$, and for all u in the model M , if $\mathcal{R}_{[i]}vu$, then $\mathcal{R}_{\otimes_i}wu$. Suppose further that $\mathcal{R}_{[i]}vz$ for an arbitrary z ; from this, and the previous statement, we may conclude that $\mathcal{R}_{\otimes_i}wz$ holds, which by the initial assumption implies that $M, z \models \phi$. Therefore, $M, v \models [i]\phi$, and since $\mathcal{R}_{\Box}wv$ holds for some v , we have that $M, w \models \Diamond [i]\phi$.

A15. Assume $M, w \models \Diamond \otimes_i \phi$. Thus, there exists a u such that $\mathcal{R}_{\Box}wu$ and $M, u \models \otimes_i \phi$. Consider an arbitrary v and z such that $\mathcal{R}_{\Box}wv$ and $\mathcal{R}_{\otimes_i}vz$. By condition (D10), and the fact that $\mathcal{R}_{\Box}wu$, $\mathcal{R}_{\Box}wv$, and $\mathcal{R}_{\otimes_i}vz$ hold, we may conclude that $\mathcal{R}_{\otimes_i}uz$ holds. Consequently, $M, z \models \phi$ holds; this fact, in conjunction with the assumption that $\mathcal{R}_{\Box}wv$ and $\mathcal{R}_{\otimes_i}vz$ hold for arbitrary v and z , implies that $M, w \models \Box \otimes_i \phi$.

A16. Assume $M, w \models \Box([i]\phi \rightarrow [i]\psi)$, $M, w \models \otimes_i \phi$, and $\mathcal{R}_{\otimes_i}wu$ for an arbitrary u . By condition (D11), the assumption $\mathcal{R}_{\otimes_i}wu$, implies that there exists a world v such that (i) $\mathcal{R}_{\Box}wv$, (ii) $\mathcal{R}_{[i]}vu$, and (iii) for all z , if $\mathcal{R}_{[i]}vz$, then $\mathcal{R}_{\otimes_i}wz$. The initial assumption, along with fact (i) that $\mathcal{R}_{\Box}wv$, entails that $M, v \models [i]\phi \rightarrow [i]\psi$. Suppose that $\mathcal{R}_{[i]}vx$ for an arbitrary x ; from fact (iii) we may conclude that $\mathcal{R}_{\otimes_i}wx$, which with the assumption that $M, w \models \otimes_i \phi$, implies

that $M, z \models \phi$. Hence, $M, v \models [i]\phi$, implying that $M, v \models [i]\psi$. Last, since we know that $\mathcal{R}_{[i]}vu$ by fact (ii), we can conclude that $M, u \models \psi$. Therefore, $M, w \models \otimes_i \phi \rightarrow \otimes_i \psi$.

Last, we show *soundness* of the **IRR**-rule from **Tstit**. Recall the rule:

$$\frac{\Box \neg p \wedge \Box (\mathbf{G}p \wedge \mathbf{H}p) \rightarrow \phi}{\phi} \text{ if } p \text{ is atomic and does not occur in } \phi$$

We assume that p does not occur in ϕ . We prove the result by contraposition and assume that ϕ is invalid. Therefore, we know there exists a model $M = (F, V)$ s.t. F is a TDS-frame and $M, w \not\models \phi$ for some $w \in W$ of M . We define another TDS-model $M' = (F, V')$ over the frame F and define the valuation V' as follows:

$$V'(q) := \begin{cases} V(q) & \text{if } q \neq p, \\ W \setminus \mathcal{R}_\Box(w) & \text{otherwise.} \end{cases}$$

where $\mathcal{R}_\Box(w) = \{v \mid (w, v) \in \mathcal{R}_\Box\}$ (i.e. the valuation V' of p contains all worlds except for those sharing the same moment with w). Clearly, since ϕ does not contain p and the other atomic propositions are valued in the same way in M as in M' we get that $M', w \models \neg \phi$. However, by the construction of V' and because F is irreflexive by condition (T7), we have that $M', w \models \Box \neg p \wedge \Box (\mathbf{G}p \wedge \mathbf{H}p)$ (the irreflexivity excludes the possibility that for some $u \in \mathcal{R}_\Box(w)$, $M', u \models p \wedge \neg p$). Since, $M', w \not\models \phi$, by Definition 3, we have that $M', w \not\models (\Box \neg p \wedge \Box (\mathbf{G}p \wedge \mathbf{H}p)) \rightarrow \phi$. Hence, we conclude that $(\Box \neg p \wedge \Box (\mathbf{G}p \wedge \mathbf{H}p)) \rightarrow \phi$ is invalid as well.

Lemma 10. *Let Γ be a MCS. Then, Γ has the following properties:*

- $\Gamma \vdash_{\text{TDS}} \phi$ iff $\phi \in \Gamma$;
- $\phi \in \Gamma$ iff $\neg \phi \notin \Gamma$;
- $\phi \wedge \psi \in \Gamma$ iff $\phi \in \Gamma$ and $\psi \in \Gamma$.

Proof. We prove each of the claims in turn:

- (i) Assume that $\phi \notin \Gamma$. Since Γ is a maximal, we know that $\Gamma \cup \{\phi\}$ is inconsistent, i.e., $\Gamma \vdash_{\text{TDS}} \phi \rightarrow \perp$. Due to the fact that Γ is consistent, we know that $\Gamma \not\vdash_{\text{TDS}} \phi$. For the opposite direction observe that if $\phi \in \Gamma$, then trivially $\Gamma \vdash_{\text{TDS}} \phi$.
- (ii) Suppose that $\phi \in \Gamma$. Observe that if $\neg \phi \in \Gamma$ as well, then Γ would be inconsistent; hence, $\neg \phi \notin \Gamma$. For the backwards direction, assume that $\neg \phi \notin \Gamma$. If $\phi \notin \Gamma$ as well, then since Γ is a MCS, we know that both $\Gamma \cup \{\phi\} \vdash_{\text{TDS}} \perp$ and $\Gamma \cup \{\neg \phi\} \vdash_{\text{TDS}} \perp$. However, this implies that $\Gamma \vdash_{\text{TDS}} \phi \wedge \neg \phi$, thus contradicting the consistency of Γ . This implies that $\phi \in \Gamma$.
- (iii) If $\phi \wedge \psi \in \Gamma$, then by fact (i) $\phi \in \Gamma$ and $\psi \in \Gamma$ since both ϕ and ψ are derivable from Γ when $\phi \wedge \psi \in \Gamma$. The opposite direction is proved similarly.

Lemma 11. *Let $\langle \alpha \rangle$ be dual to $[\alpha] \in \text{Boxes}$. Then, $\mathcal{R}_{[\alpha]} \Gamma \Delta$ iff for all $\phi \in \mathcal{L}_{\text{TDS}}$, if $\phi \in \Delta$, then $\langle \alpha \rangle \phi \in \Gamma$.*

Proof. Let $\langle \alpha \rangle$ be dual to $[\alpha] \in \text{Boxes}$ and let Γ and Δ be maximally consistent IRR-theories. We prove both directions of the equivalence.

First, assume that $\mathcal{R}_{[\alpha]} \Gamma \Delta$ holds and consider an arbitrary $\phi \in \Delta$. Since Δ is a MCS, we know that $\neg \phi \notin \Delta$, which implies by the definition of $\mathcal{R}_{[\alpha]}$ that $[\alpha] \neg \phi \notin \Gamma$. Due to the fact that Γ is a MCS, this implies that $\neg[\alpha] \neg \phi \in \Gamma$, which further implies that $\langle \alpha \rangle \phi \in \Gamma$.

For the opposite direction of the equivalence assume that for all $\phi \in \mathcal{L}_{\text{TDS}}$, if $\phi \in \Delta$, then $\langle \alpha \rangle \phi \in \Gamma$. Let $\psi \in \mathcal{L}_{\text{TDS}}$ and assume that $[\alpha] \psi \in \Gamma$. Then, since Γ is a MCS, we know that $\langle \alpha \rangle \neg \psi \notin \Gamma$. Therefore, $\neg \psi \notin \Delta$, which implies that $\psi \in \Delta$ since Δ is a MCS. Since ψ was arbitrary, we have established that $\mathcal{R}_{[\alpha]} \Gamma \Delta$.

Lemma 2 *Let $\phi \in \mathcal{L}_{\text{TDS}}$ be a consistent formula. Then, there exists an IRR-theory Γ such that $\phi \in \Gamma$.*

Proof. Let $\phi \in \mathcal{L}_{\text{TDS}}$ be a consistent formula. We enumerate the formulae of \mathcal{L}_{TDS} so that each formula in odd position is an element of **Zig** and make use of this enumeration to build an increasing sequence of consistent theories $\Gamma_0, \Gamma_1, \dots, \Gamma_n, \dots$

We let $\Gamma_0 := \{\phi \wedge \Box \neg p \wedge \Box(\text{G}p \wedge \text{H}p)\}$ for some propositional variable p not occurring in ϕ . We define the sequence of Γ_n (for $n > 0$) as follows: Assume that Γ_n is defined and consider ψ_n of the enumeration. We know that either $\Gamma_n \cup \{\neg \psi_n\}$ is consistent or $\Gamma_n \cup \{\psi_n\}$ is consistent. If $\Gamma_n \cup \{\neg \psi_n\}$ is consistent, set $\Gamma_{n+1} := \Gamma_n \cup \{\neg \psi_n\}$. If $\Gamma_n \cup \{\psi_n\}$ is consistent, then there are two cases to consider: either (i) n is even or (ii) n is odd. If n is even, then set $\Gamma_{n+1} := \Gamma_n \cup \{\psi_n\}$. Otherwise, set $\Gamma_{n+1} := \Gamma_n \cup \{\psi_n, \psi_n(q)\}$, where q is a propositional variable not occurring in Γ_n or ψ . We define our desired maximally consistent IRR-theory as follows:

$$\Gamma := \bigcup_{n \in \mathbb{N}} \Gamma_n$$

To finish the proof we need to show that Γ is both a MCS and IRR-theory. We first prove that (i) Γ is a MCS and then show that (ii) Γ is an IRR-theory.

To prove claim (i), it is useful to first prove that for all $n \in \mathbb{N}$, each Γ_n is consistent. We show this claim by induction on n . In the base case, assume for a contradiction that $\Gamma_0 = \{\phi \wedge \Box \neg p \wedge \Box(\text{G}p \wedge \text{H}p)\}$ is inconsistent. Hence, $\Box \neg p \wedge \Box(\text{G}p \wedge \text{H}p) \wedge \phi \vdash_{\text{TDS}} \perp$, which further implies that $\vdash_{\text{TDS}} \Box \neg p \wedge \Box(\text{G}p \wedge \text{H}p) \rightarrow (\phi \rightarrow \perp)$. We may infer from the rule R2 that $\vdash_{\text{TDS}} \phi \rightarrow \perp$. However, we know that ϕ is consistent, meaning that $\not\vdash_{\text{TDS}} \phi \rightarrow \perp$. We have thus obtained a contradiction implying then that Γ_0 is in fact consistent. For the inductive step assume that Γ_n is consistent. We want to show that Γ_{n+1} is consistent. This trivially follows by the definition of Γ_{n+1} .

To prove that Γ is a MCS, we must show that Γ is both consistent and maximal. Assume for a contradiction that Γ is inconsistent. Then, this implies that for some finite subset Γ' of Γ , $\Gamma' \vdash \perp$. However, if this is the case, then there exists some Γ_n such that $\Gamma_n \vdash_{\text{TDS}} \perp$. We know that this cannot be the case by the previous paragraph, and so, Γ must be consistent. Assume now that there exists some Γ' such that $\Gamma \subset \Gamma'$ and $\Gamma' \not\vdash_{\text{TDS}} \perp$. Let $\psi \in \Gamma' \setminus \Gamma$. Since ψ

is a formula in \mathcal{L}_{TDS} , we know that if was considered at some point during the construction of the sequence $\Gamma_0, \Gamma_1, \dots, \Gamma_n, \dots$. Since $\psi \notin \Gamma$ this implies that there exists some Γ_m such that $\Gamma_m \cup \{\psi\}$ is inconsistent. Therefore, $\Gamma_m \vdash_{\text{TDS}} \neg\psi$, which implies that $\Gamma \vdash_{\text{TDS}} \neg\psi$. Due to the fact that $\Gamma \subset \Gamma'$, it follows that $\Gamma' \vdash_{\text{TDS}} \neg\psi$ and $\Gamma' \vdash_{\text{TDS}} \psi$ since $\psi \in \Gamma'$, which is a contradiction. Therefore, Γ is a MCS.

We now prove that Γ is an IRR-theory. By construction we know that $\phi \wedge \Box\neg p \wedge \Box(\mathsf{G}p \wedge \mathsf{H}p) \in \Gamma_0 \subset \Gamma$, and since Γ is a MCS, it follows that $\Box\neg p \wedge \Box(\mathsf{G}p \wedge \mathsf{H}p) \in \Gamma$, thus satisfying the first condition of being an IRR-theory. The second condition of being an IRR-theory is satisfied by the fact that whenever a formula $\psi \in \text{Zig}$ is added to $\Gamma_m \subset \Gamma$, for $m \in \mathbb{N}$, the formula $\psi(q)$ is added as well with q fresh.

Lemma 3 *Let Γ be an IRR-theory and let $\langle \alpha \rangle$ be dual to $[\alpha] \in \text{Boxes}$. For each $\langle \alpha \rangle \phi \in \Gamma$ there exists an IRR-theory Δ such that $\mathcal{R}_{[\alpha]} \Gamma \Delta$.*

Proof. Similar to [14, Lem. 16].

Lemma 6 *Let Γ be an IRR-theory in W . Then, there exists an IRR-theory $\Delta \in W$ such that $\mathcal{R}_{\Box}^{dt} \Gamma \Delta$ and for every IRR-theory $\Sigma \in W^{dt}$, if $\mathcal{R}_{[i]}^{dt} \Delta \Sigma$, then $\mathcal{R}_{\otimes_i}^{dt} \Gamma \Sigma$.*

Proof. Let Γ be an arbitrary IRR-theory in W^{dt} . Since Γ is an IRR-theory, there is a propositional variable p such that $\text{name}(p) \in \Gamma$. Define

$$\Delta_0 := \{[i]\phi \mid \otimes_i \phi \in \Gamma\} \cup \{\psi \mid \Box\psi \in \Gamma\} \cup \{\text{name}(p)\}.$$

We will prove by contradiction that Δ_0 is consistent and then extend Δ_0 to an IRR-theory.

If Δ_0 is inconsistent, then

$$\vdash_{\text{TDS}} ([i]\phi_1 \wedge \dots \wedge [i]\phi_n \wedge \psi_1 \wedge \dots \wedge \psi_n \wedge \text{name}(p)) \rightarrow \perp$$

where $\psi_1, \dots, \psi_n \in \{\psi \mid \Box\psi \in \Gamma\}$ and $[i]\phi_1, \dots, [i]\phi_n \in \{[i]\phi \mid \otimes_i \phi \in \Gamma\}$. Let $\hat{\phi} = \phi_1 \wedge \dots \wedge \phi_n$ and $\hat{\psi} = \psi_1 \wedge \dots \wedge \psi_n$. Since, $\vdash_{\text{TDS}} [i]\hat{\phi} \leftrightarrow [i]\phi_1 \wedge \dots \wedge [i]\phi_n$ we get

$$\vdash_{\text{TDS}} \hat{\psi} \wedge \text{name}(p) \rightarrow \neg[i]\hat{\phi}$$

By necessitation for \Box and the \Box K-axiom, we get $\vdash_{\text{TDS}} \Box(\hat{\psi} \wedge \text{name}(p)) \rightarrow \Box\neg[i]\hat{\phi}$, which implies $\vdash_{\text{TDS}} \Box\hat{\psi} \wedge \Box\text{name}(p) \rightarrow \neg\Diamond[i]\hat{\phi}$. Clearly, because $\Box\hat{\psi} \in \Gamma$, $\text{name}(p) \in \Gamma$ and $\vdash_{\text{TDS}} \text{name}(p) \rightarrow \Box\text{name}(p)$, we have that $\Gamma \vdash_{\text{TDS}} \neg\Diamond[i]\hat{\phi}$. This implies that $\neg\Diamond[i]\hat{\phi} \in \Gamma$ since Γ is an IRR-theory.

Also, since $\otimes_i \phi_1, \dots, \otimes_i \phi_n \in \Gamma$ we have $\otimes_i \phi_1 \wedge \dots \wedge \otimes_i \phi_n \in \Gamma$ since Γ is an IRR-theory. By $\vdash_{\text{TDS}} \otimes_i \hat{\phi} \leftrightarrow \otimes_i \phi_1 \wedge \dots \wedge \otimes_i \phi_n$ we conclude $\otimes_i \hat{\phi} \in \Gamma$ as well. Since $\otimes_i \hat{\phi} \rightarrow \Diamond[i]\hat{\phi} \in \Gamma$ because the formula is an instance of axiom A14, we obtain by modus ponens that $\Diamond[i]\hat{\phi} \in \Gamma$. Since Γ is an IRR-theory (and hence consistent) we obtain a contradiction, which proves that Δ_0 is consistent.

We now extend Δ_0 to an IRR-theory Δ by first defining an increasing sequence $\Delta_0, \Delta_1, \dots, \Delta_n, \dots$ of sets of formulae. Suppose that Δ_n is consistent

and defined, and enumerate the formulae of \mathcal{L}_{TDS} so that each formula in odd position is an element of **Zig**; we aim to define Δ_{n+1} .

Consider the formula ψ_n . Either, $\Delta_n \cup \{\neg\psi_n\}$ is consistent or $\Delta_n \cup \{\psi_n\}$ is consistent. If the former holds, then set $\Delta_{n+1} := \Delta_n \cup \{\neg\psi_n\}$. If the latter holds, then there are two subcases to consider: either n is even, in which case, we set $\Delta_{n+1} := \Delta_n \cup \{\psi_n\}$, or n is odd, in which case, $\Delta_n \cup \{\psi_n\}$ is consistent and $\psi_n \in \mathbf{Zig}$. We show that in the latter subcase we can find a propositional variable q such that $\Delta_n \cup \{\psi_n, \psi_n(q)\}$ is consistent; we then define $\Delta_{n+1} := \Delta_n \cup \{\psi_n, \psi_n(q)\}$.

Observe that

$$\ominus_i (\text{name}(p) \wedge \bigwedge_{\chi \in \Delta_n \setminus \Delta_0} \chi \wedge \psi_n) \in \Gamma \quad (3)$$

For otherwise,

$$\otimes_i ((\text{name}(p) \wedge \bigwedge_{\chi \in \Delta_n \setminus \Delta_0} \chi) \rightarrow \neg\psi_n) \in \Gamma$$

since Γ is an IRR-theory and has the properties specified by Lemma 10. By the definition of Δ_0 it follows that

$$[i]((\text{name}(p) \wedge \bigwedge_{\chi \in \Delta_n \setminus \Delta_0} \chi) \rightarrow \neg\psi_n) \in \Delta_n$$

Using the fact that $\vdash_{\text{TDS}} [i]\theta \rightarrow \theta$ holds for any formula θ , we infer that

$$\Delta_n \vdash_{\text{TDS}} (\text{name}(p) \wedge \bigwedge_{\chi \in \Delta_n \setminus \Delta_0} \chi) \rightarrow \neg\psi_n$$

Since

$$\Delta_n \vdash_{\text{TDS}} \text{name}(p) \wedge \bigwedge_{\chi \in \Delta_n \setminus \Delta_0} \chi$$

we may conclude that $\Delta_n \vdash_{\text{TDS}} \neg\psi_n$, which contradicts the fact that $\Delta_n \cup \{\psi_n\}$ is consistent. Therefore, since Γ is an IRR-theory and (1) holds, we know that

$$\ominus_i (\text{name}(p) \wedge \bigwedge_{\chi \in \Delta_n \setminus \Delta_0} \chi \wedge \psi_n(q)) \in \Gamma \quad (4)$$

Using this fact, we may prove that $\Delta_{n+1} := \Delta_n \cup \{\psi_n, \psi_n(q)\}$ is consistent, for suppose otherwise. Then, there exist $\zeta_1, \dots, \zeta_m \in \{\zeta \mid \Box\zeta \in \Gamma\}$ and $[i]\xi_1, \dots, [i]\xi_k \in \{[i]\xi \mid \otimes_i \xi \in \Gamma\}$ such that

$$\vdash_{\text{TDS}} \zeta_1 \wedge \dots \wedge \zeta_m \rightarrow ([i]\xi_1 \wedge \dots \wedge [i]\xi_k \rightarrow \neg(\text{name}(p) \wedge \bigwedge_{\chi \in \Delta_n \setminus \Delta_0} \chi \wedge \psi_n(q)))$$

By \otimes_i necessitation and the \otimes_i K-axiom, we can derive

$$\vdash_{\text{TDS}} \otimes_i (\zeta_1 \wedge \dots \wedge \zeta_m) \rightarrow \otimes_i ([i]\xi_1 \wedge \dots \wedge [i]\xi_k \rightarrow \neg(\text{name}(p) \wedge \bigwedge_{\chi \in \Delta_n \setminus \Delta_0} \chi \wedge \psi_n(q)))$$

Using axiom A13 we obtain

$$\vdash_{\text{TDS}} \Box(\zeta_1 \wedge \cdots \wedge \zeta_m) \rightarrow \otimes_i([i]\xi_1 \wedge \cdots \wedge [i]\xi_k \rightarrow \neg(\text{name}(p) \wedge \bigwedge_{\chi \in \Delta_n \setminus \Delta_0} \chi \wedge \psi_n(q)))$$

By our assumption and the fact that Γ is an IRR-theory, we know that $\Box(\zeta_1 \wedge \cdots \wedge \zeta_m) \in \Gamma$, implying that

$$\otimes_i([i]\xi_1 \wedge \cdots \wedge [i]\xi_k \rightarrow \neg(\text{name}(p) \wedge \bigwedge_{\chi \in \Delta_n \setminus \Delta_0} \chi \wedge \psi_n(q))) \in \Gamma$$

We infer the following using modal reasoning

$$\otimes_i[i](\xi_1 \wedge \cdots \wedge \xi_k) \rightarrow \otimes_i \neg(\text{name}(p) \wedge \bigwedge_{\chi \in \Delta_n \setminus \Delta_0} \chi \wedge \psi_n(q))) \in \Gamma$$

One can confirm that $\vdash_{\text{TDS}} \otimes_i \theta \rightarrow \otimes_i[i]\theta$ (See [15]) holds for any formula θ , and therefore

$$\otimes_i(\xi_1 \wedge \cdots \wedge \xi_k) \rightarrow \otimes_i \neg(\text{name}(p) \wedge \bigwedge_{\chi \in \Delta_n \setminus \Delta_0} \chi \wedge \psi_n(q))) \in \Gamma$$

Our assumption implies that $\otimes_i(\xi_1 \wedge \cdots \wedge \xi_k) \in \Gamma$, and so

$$\otimes_i \neg(\text{name}(p) \wedge \bigwedge_{\chi \in \Delta_n \setminus \Delta_0} \chi \wedge \psi_n(q))) \in \Gamma$$

This contradicts (2) and proves that $\Delta_n \cup \{\psi_n \psi_n(q)\}$ is consistent.

It is easy to infer that Δ is an IRR-theory by an argument similar to Lemma 2.

Clearly, $\mathcal{R}_{\Box}^{dt} \Gamma \Delta$ holds by the definition of Δ . Last, let Σ be an arbitrary IRR-theory in W^{dt} . Assume that $\mathcal{R}_{[i]}^{dt} \Delta \Sigma$ holds and let $\otimes_i \xi \in \Gamma$. By definition $[i]\xi \in \Delta$, and so, $\xi \in \Sigma$ by the definition of the relation $\mathcal{R}_{[i]}^{dt}$, which completes the proof.

Lemma 7 *Let Γ and Δ be IRR-theories in W^{dt} such that $\mathcal{R}_{\otimes_i}^{dt} \Gamma \Delta$. Then, there exists an IRR-theory $\Sigma \in W$ such that $\mathcal{R}_{\Box}^{dt} \Gamma \Sigma$, $\mathcal{R}_{[i]}^{dt} \Sigma \Delta$, and for all $\Pi \in W^{dt}$, if $\mathcal{R}_{[i]}^{dt} \Sigma \Pi$, then $\mathcal{R}_{\otimes_i}^{dt} \Gamma \Pi$.*

Proof. To prove this lemma, we proceed differently compared to Lemma 6, making explicit use of the existence lemma (Lemma 3). Let Γ and Δ be IRR-theories in W^{dt} such that $\mathcal{R}_{\otimes_i}^{dt} \Gamma \Delta$. Then, there is a $\text{name}(p)$ for some p such that $\text{name}(p) \in \Delta$. Since $\phi \rightarrow \langle i \rangle \phi \in \Delta$ for any $\phi \in \mathcal{L}_{\text{TDS}}$ we know $\langle i \rangle \text{name}(p) \in \Delta$. Hence, by Lemma 3 we know there exists a $\Sigma \in W^{dt}$ for which $\mathcal{R}_{[i]}^{dt} \Delta \Sigma$. First, we show (i) $\mathcal{R}_{[i]}^{dt} \Sigma \Delta$, then we show (ii) $\mathcal{R}_{\Box}^{dt} \Gamma \Sigma$ and last we show (iii) for any $\Pi \in W^{dt}$ for which $\mathcal{R}_{[i]}^{dt} \Sigma \Pi$, we have $\mathcal{R}_{\otimes_i}^{dt} \Gamma \Pi$.

- (i) Recall $\mathcal{R}_{[i]}^{dt} \Delta \Sigma$, take an arbitrary $[i]\phi \in \Sigma$, it suffices to show that $\phi \in \Delta$. By Lemma 11, we know that $\langle i \rangle [i]\phi \in \Delta$. Since $\vdash_{\text{TDS}} \langle i \rangle [i]\theta \rightarrow \theta$ for any $\theta \in \mathcal{L}_{\text{TDS}}$ (by axiom A5, A6, and propositional reasoning) we obtain $\phi \in \Delta$; hence $\mathcal{R}_{[i]}^{dt} \Sigma \Delta$.
- (ii) Assume an arbitrary $\Box\phi \in \Gamma$. We prove that $\phi \in \Sigma$. We know $\vdash_{\text{TDS}} \Box\phi \rightarrow \otimes_i \phi$ (axiom A13). Hence, since Γ is an IRR-theory, we obtain $\otimes_i \phi \in \Gamma$. Furthermore, $\vdash_{\text{TDS}} \otimes_i \phi \rightarrow \otimes_i [i]\phi$ (See [15]), and therefore, $\otimes_i [i]\phi \in \Gamma$. Since $\mathcal{R}_{\otimes_i}^{dt} \Gamma \Delta$ we get $[i]\phi \in \Delta$ and thus, by the fact that $\mathcal{R}_{[i]}^{dt} \Delta \Sigma$, we know $\phi \in \Sigma$. We conclude $\mathcal{R}_{\Box}^{dt} \Gamma \Sigma$.
- (iii) Take an arbitrary $\Pi \in W^{dt}$. Assume $\mathcal{R}_{[i]}^{dt} \Sigma \Pi$. and $\otimes_i \phi \in \Gamma$. Since $\otimes_i \phi \rightarrow \otimes_i [i]\phi \in \Gamma$, $\otimes_i [i]\phi \in \Gamma$. Furthermore, since $\vdash_{\text{TDS}} \otimes_i [i]\theta \rightarrow \otimes_i [i][i]\theta$ for any $\theta \in \mathcal{L}_{\text{TDS}}$ (A5, A6, R1, A12), we know $\otimes_i [i][i]\phi \in \Gamma$, and thus $[i][i]\phi \in \Delta$. Consequently, we get $[i]\phi \in \Sigma$ and last $\phi \in \Pi$, giving us $\mathcal{R}_{\otimes_i}^{dt} \Gamma \Pi$.

Lemma 8 *The canonical model $M^{dt}|_{\text{IRR}}$ belongs to the class of TDS models.*

Proof. The argument that $M^{dt}|_{\text{IRR}}$ possesses properties (C1), (C2), (C3)*, (T4)–(T7) is the same as in [14, Lem. 19]. Therefore, we need only confirm that the model satisfies conditions (D8)–(D11).

The fact that $M^{dt}|_{\text{IRR}}$ satisfies conditions (D9) and (D11) follows from Lemma 6 and 7. We additionally prove that $M^{dt}|_{\text{IRR}}$ satisfies conditions (D8) and (D10).

(D8) Let Γ and Δ be arbitrary IRR-theories. Assume that $\mathcal{R}_{\otimes_i}^{dt} \Gamma \Delta$ and assume that $\phi \in \Delta$. Hence, by Lemma 11, we know that $\ominus_i \phi \in \Gamma$. Since $\Box \neg \phi \rightarrow \otimes_i \neg \phi \in \Gamma$, we have $\ominus_i \phi \rightarrow \Diamond \phi \in \Gamma$. Hence, $\Diamond \phi \in \Gamma$, which implies that $\mathcal{R}_{\Box}^{dt} \Gamma \Delta$.

(D10) Let $\Gamma, \Delta, \Sigma, \Pi \in W^{dt} \cap \text{IRR}$ and assume that $\mathcal{R}_{\Box}^{dt} \Gamma \Delta$, $\mathcal{R}_{\Box}^{dt} \Gamma \Sigma$, and $\mathcal{R}_{\otimes_i}^{dt} \Sigma \Pi$. We will show that $\mathcal{R}_{\otimes_i}^{dt} \Delta \Pi$.

Let $\phi \in \mathcal{L}_{\text{TDS}}$ and assume $\phi \in \Pi$. Then $\ominus_i \phi \in \Sigma$ and, hence, $\Diamond \ominus_i \phi \in \Gamma$ by Lemma 11. Since

$$\vdash_{\text{TDS}} (\Diamond \otimes_i \phi \rightarrow \Box \otimes_i \phi) \rightarrow (\Diamond \ominus_i \phi \rightarrow \Box \ominus_i \phi)$$

and

$$\Diamond \otimes_i \phi \rightarrow \Box \otimes_i \phi \in \Gamma$$

we may infer that $\Diamond \ominus_i \phi \rightarrow \Box \ominus_i \phi \in \Gamma$. Due to the fact that $\Diamond \ominus_i \phi \in \Gamma$, we obtain $\Box \ominus_i \phi \in \Gamma$, and so, $\ominus_i \phi \in \Delta$. Therefore, $\mathcal{R}_{\otimes_i}^{dt} (\Delta, \Pi)$.

Theorem 2 *If $\phi \in \mathcal{L}_{\text{TDS}}$ is a consistent formula, then ϕ is satisfiable on a TDS-model.*

Proof. Suppose that $\phi \in \mathcal{L}_{\text{TDS}}$ is consistent. By Lemma 2, we can extend ϕ to an IRR-theory Γ such that $\phi \in \Gamma$. By Lemma 3, we know that the set IRR is a diamond saturated set, and so, by Lemma 1, we know that $M^{dt}|_{\text{IRR}}, \Gamma \models \phi$ iff $\phi \in \Gamma$. Hence, we can conclude that $M^{dt}|_{\text{IRR}}, \Gamma \models \phi$. By Lemma 8 we know that $M^{dt}|_{\text{IRR}}$ is a TDS-model; therefore, ϕ is satisfiable on a TDS-model.