

# Moorean Phenomena in Epistemic Logic<sup>\*</sup>

Wesley H. Holliday and Thomas F. Icard, III

*Department of Philosophy  
Stanford University  
Stanford, California, USA*

---

## Abstract

A well-known open problem in epistemic logic is to give a syntactic characterization of the *successful* formulas. Semantically, a formula is successful if and only if for any pointed model where it is true, it remains true after deleting all points where the formula was false. The classic example of a formula that is not successful in this sense is the “Moore sentence”  $p \wedge \neg \Box p$ , read as “ $p$  is true but you do not know  $p$ .” Not only is the Moore sentence unsuccessful, it is *self-refuting*, for it never remains true as described. We show that in logics of knowledge and belief for a single agent (extended by **S5**), Moorean phenomena are the source of all self-refutation; moreover, in logics for an *introspective* agent (extending **KD45**), Moorean phenomena are the source of all unsuccessfulness as well. This is a distinctive feature of such logics, for with a non-introspective agent or multiple agents, non-Moorean unsuccessful formulas appear. We also consider how successful and self-refuting formulas relate to the *Cartesian* and *learnable* formulas, which have been discussed in connection with Fitch’s “paradox of knowability.” We show that the Cartesian formulas are exactly the formulas that are not *eventually* self-refuting and that not all learnable formulas are successful. In an appendix, we give syntactic characterizations of the successful and the self-refuting formulas.

*Keywords:* epistemic logic, successful formulas, Moore sentence

---

## 1 Introduction

According to the epistemic interpretation of modal logic, the points in a modal model represent ways the world might be, consistent with an agent’s information. In this context, “learning” a formula amounts to eliminating those points in the model where the formula is false. The resulting submodel represents the agent’s information state after learning has occurred. Some formulas—though not all—remain true whenever they are learned. A well-known open problem [12,4,5,6,7,2] in epistemic logic is to give a syntactic characterization of these *successful* formulas. Partial results have been obtained (Section 2), but a full solution has proven elusive.

The classic example of an *unsuccessful* formula is the Moore sentence  $p \wedge \neg \Box p$ , read as “ $p$  is true but you do not know  $p$ .” This example is a second-person variation of G.E. Moore’s famous puzzle [16] involving the paradoxical first-person assertion, “ $p$  is true but I do not believe  $p$ .” Hintikka devoted a chapter of his seminal 1962 monograph *Knowledge and Belief* [13] to an analysis of such sentences, including the second-person Moore sentence. Hintikka observed its unsuccessfulness as follows: “You may come to know that what I said *was* true, but saying it in so many words

---

<sup>\*</sup> In L. Beklemishev, V. Goranko and V. Shehtman, eds., *Advances in Modal Logic*, Volume 8, 178-199, College Publications, 2010. Minor changes to the published version have been made in this version.

has the effect of making what is being said false” [13, p. 69]. Yet the formal question of unsuccessfulness did not arise for Hintikka. Only with the advent of Dynamic Epistemic Logic (see, e.g., [7]) and the idea of learning as model reduction have the Moorean phenomena and unsuccessfulness been formally related.

The reason the Moore sentence is unsuccessful is that it can only be true in a model if  $p$  is true at some point and false at some other point accessible from the first. When the agent learns the sentence, all points where  $p$  is false are eliminated from the model, including all witnesses for  $\neg\Box p$ , so the sentence becomes false. This shows that the Moore sentence is not only unsuccessful, it is *self-refuting*, for it always becomes false when learned. Related to the self-refuting property of the Moore sentence is the fact that the sentence cannot be known. Indeed, the Moore sentence is at the root of Fitch’s famous “paradox of knowability” [11,3,9]: if there is an unknown truth, then there is an unknowable truth. For if  $p$  is true but unknown, then the Moore sentence  $p \wedge \neg\Box p$  is true, but the Moore sentence cannot be known, because  $\Box(p \wedge \neg\Box p)$  is inconsistent with standard assumptions about knowledge.

While the Moore sentence is conspicuously unsuccessful, other unsuccessful formulas are less conspicuous. The formula  $\neg(p \vee q) \vee (p \wedge (\Box p \vee \Diamond q))$  is also unsuccessful, as we show in Example 5.12 below, but is the reason Moorean? While on the surface this formula looks unlike a Moore sentence, it is in fact possible to transform the formula to reveal its Moorean character. Indeed, we will prove that for a wide range of logics, such a transformation is possible for every unsuccessful formula. However well-disguised, their nature is always Moorean.

In Section 2 we establish notation, give the definitions of *successful*, *self-refuting*, etc., and review what is already known in the literature on successful formulas. In Section 3 we show that in logics of knowledge and belief for a single agent (extended by **S5**), Moorean phenomena are the source of all self-refutation; moreover, in logics for an *introspective* agent (extending **KD45**), Moorean phenomena are the source of all unsuccessfulness as well. This is a distinctive feature of such logics, for as we show in Section 4, in logics for a non-introspective agent or multiple agents, non-Moorean unsuccessful formulas appear. Finally, in Section 5 we relate successful and self-refuting formulas to the *Cartesian* and *learnable* formulas, which have been discussed in connection with Fitch’s paradox, and to the *informative*, *eventually self-refuting*, *super-successful* formulas, which we introduce here. In Appendix A we give syntactic characterizations of the successful and the self-refuting formulas.

## 2 Preliminaries and Previous Results

Throughout we work with a fixed, unimodal language with  $\Box$  and its dual  $\Diamond$ , and an infinite set  $\mathbf{Prop}$  of propositional variables. The expression  $\Diamond^+\varphi$  abbreviates  $\Diamond\varphi \wedge \varphi$ . We use the standard semantics of modal logic, where a model is a triple  $\langle W, R, V \rangle$  with  $W$  any set of points,  $R \subseteq W \times W$  any relation, and  $V : \mathbf{Prop} \rightarrow \wp(W)$  a valuation function. A pointed model is a pair  $\mathcal{M}, w$  with  $\mathcal{M}$  a model and  $w \in W$ . The satisfaction relation  $\models$  between pointed models and formulas is defined as usual.

It is typical to take **KD45** as a logic of belief and **S5** as a logic of knowledge. Apart from the assumptions that one does not believe or know anything inconsistent, and that what is known is true, these logics assume *positive* (axiom 4) and

*negative* (axiom **5**) *introspection*: if one believes something, then one believes that one believes it; if one does not believe something, then one believes that one does not believe it; and *mutatis mutandis* for knowledge. Most work on successful formulas has assumed **S5**. Since our results apply to a wider range of logics, we assume through Section 3 that we are working with at least **KD45**, so all of our models will be serial, transitive, and Euclidean. We call such models *quasi-partitions*.

Where  $\mathcal{M}'$  is a submodel (not necessarily proper) of  $\mathcal{M}$ , we write  $\mathcal{M}' \subseteq \mathcal{M}$ . Rather than studying formulas preserved under arbitrary submodels, we will study formulas preserved under a special way of taking submodels, as in the following.

**Definition 2.1** Given a model  $\mathcal{M} = \langle W, R, V \rangle$ , the *relativization of  $\mathcal{M}$  to  $\varphi$*  is the (possibly empty) submodel  $\mathcal{M}_{|\varphi} = \langle W_{|\varphi}, R_{|\varphi}, V_{|\varphi} \rangle$  of  $\mathcal{M}$ , where  $W_{|\varphi} = \{w \in W : \mathcal{M}, w \models \varphi\}$ ,  $R_{|\varphi}$  is  $R$  restricted to  $W_{|\varphi}$ , and  $V_{|\varphi}(p) = V(p) \cap W_{|\varphi}$ .

**Definition 2.2** A formula  $\varphi$  is *successful* (in logic  $\mathbf{L}$ ) iff for every pointed model (of  $\mathbf{L}$ ),  $\mathcal{M}, w \models \diamond^+ \varphi$  implies  $\mathcal{M}_{|\varphi}, w \models \varphi$ . A formula is *unsuccessful* (in  $\mathbf{L}$ ) iff it is not successful. A formula is *self-refuting* (in  $\mathbf{L}$ ) iff for every pointed model (of  $\mathbf{L}$ ),  $\mathcal{M}, w \models \diamond^+ \varphi$  implies  $\mathcal{M}_{|\varphi}, w \not\models \varphi$ .<sup>1</sup>

In the standard definitions of successful and self-refuting formulas, where  $\mathbf{L}$  is assumed to be **S5**, the precondition only requires that  $\varphi$  be true at  $w$ . Since we are also working with **KD45**, we additionally require that  $\varphi$  be true at an accessible point, so that  $\mathcal{M}_{|\varphi}$  is a quasi-partition provided  $\mathcal{M}$  is. Our definition reduces to the standard one in the case of **S5**. In either case, unsatisfiable formulas are self-refuting and successful. In the case of **KD45**, a satisfiable formula such as  $p \wedge \Box \neg p$ , read as “ $p$  is true but you believe  $\neg p$ ,” is self-refuting and successful, since  $\diamond^+ (p \wedge \Box \neg p)$  is unsatisfiable. While it may be more intuitive to require the satisfiability of  $\diamond^+ \varphi$  for a successful  $\varphi$ , we will follow the standard definition in not requiring satisfiability.

The following lemma relates success and self-refutation across different logics.

**Lemma 2.3** *Let  $\mathbf{L}$  be a sublogic of  $\mathbf{L}'$ . If  $\varphi$  is unsuccessful in  $\mathbf{L}'$ , then  $\varphi$  is unsuccessful in  $\mathbf{L}$ , and if  $\varphi$  is self-refuting in  $\mathbf{L}$ , then  $\varphi$  is self-refuting in  $\mathbf{L}'$ .*

**Proof.** Immediate from Definition 2.2, given that models of  $\mathbf{L}'$  are models of  $\mathbf{L}$ .  $\square$

An obstacle to giving a simple syntactic characterization of the set of successful formulas is its lack of closure properties. Successful formulas are not closed under negation (take  $\neg p \vee \Box p$ ), conjunction (take  $p$  and  $\neg \Box p$ ), or implication (use  $\varphi \rightarrow \perp$ ) [6]. We show in Proposition 5.11 that they are also not closed under disjunction. Conversely, if a negated formula is successful, the unnegated formula may be unsuccessful (take  $\neg(p \wedge \diamond \neg p)$ ), if a conjunction is successful, some of the conjuncts may be unsuccessful (take  $(p \wedge \diamond \neg p) \wedge \neg p$ ), and if a disjunction is successful, some or even all of the disjuncts may be unsuccessful (Proposition 5.9 and Example 5.2).

By contrast, the formulas preserved under arbitrary submodels are well-behaved. The following result was proved independently by van Benthem and Visser [20,8]. A formula is *universal* iff it can be constructed using only literals,  $\wedge$ ,  $\vee$ , and  $\Box$ .

<sup>1</sup> The term ‘successful’ is used by Gerbrandy [12], by analogy with the success postulate of belief revision, while the term ‘self-refuting’ is used by van Benthem [3]. Self-refuting formulas have also been called *strongly unsuccessful* [2].

**Theorem 2.4** *A formula is preserved under submodels (of all relational models) iff it is equivalent (in  $\mathbf{K}$ ) to a universal formula.*

Similarly, a formula is preserved under model extensions iff it is equivalent to an *existential* formula, constructed using only literals,  $\wedge$ ,  $\vee$ , and  $\diamond$  [8].

Lemma 2.3 and the right-to-left direction of Theorem 2.4 give the following.

**Corollary 2.5** *Universal formulas are successful in any normal modal logic.*

As noted by van Benthem [4], for any model  $\mathcal{M}$  and formula  $\varphi$ , there is a universal formula  $\varphi'$  such that  $\mathcal{M}|_{\varphi} = \mathcal{M}|_{\varphi'}$ .<sup>2</sup> However, this result assumes the relation for  $\mathcal{M}$  is at least a quasi-partition. For example, given the model in Figure 1, where the relation is not Euclidean, there is no universal formula  $\psi$  such that  $\mathcal{M}|_{\diamond p} = \mathcal{M}|_{\psi}$ . This is symptomatic of the fact, established in Section 4, that in logics without both axioms 4 and 5, there are non-Moorean sources of unsuccessfulness.

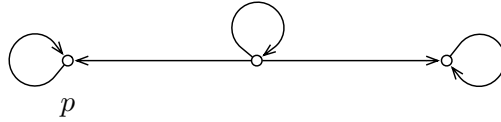


Fig. 1.

Gerbrandy [12] proved a proposition similar to Corollary 2.5, but in a slightly different formal context and with  $\diamond$  only as a defined operator. He showed that if all instances of  $\Box$  are within the scope of an even number of negations, then the formula is successful. However, the converse of Corollary 2.5 does not hold, even in the case of formulas containing no nested modal operators. For example, the formula  $\diamond p \vee p$  is successful in  $\mathbf{K}$ , but it is not equivalent to any universal formula. Yet other connections with universal formulas do hold. Qian [18] showed that a formula containing no nested modal operators is successful in  $\mathbf{K}$  iff it can be transformed using a certain algorithm into a universal formula; moreover, a “homogeneous” formula, containing no nested modal operators and no propositional variables outside the scope of modal operators, is successful in  $\mathbf{K}$  iff it is equivalent to a universal formula. The restriction to  $\mathbf{K}$  is essential, for the formula  $\diamond p$  is successful in extensions of **KD45** (not in  $\mathbf{K}$ ), but it is not equivalent to any universal formula.

As far as we know, there are no other published results on the syntax of successful formulas in the basic modal language. However, successful formulas are often studied in the more general context of Dynamic Epistemic Logic (DEL) [1,12,7], where model changing operations are formalized in the logic itself by adding dynamic operators to the language. In the simplest fragment of DEL extending **S5**, known as Public Announcement Logic (PAL), there are sentences of the form  $[\varphi]\psi$ , read as “after the *announcement* of  $\varphi$ ,  $\psi$  is true,” for which satisfaction is defined:

$$\mathcal{M}, w \models [\varphi]\psi \text{ iff } \mathcal{M}, w \models \varphi \text{ implies } \mathcal{M}|_{\varphi}, w \models \psi.$$

An advantage of the PAL setting is that it allows for simple definitions of success and self-refutation [6]. Successful formulas are those for which  $[\varphi]\varphi$  is valid (in **S5**),

<sup>2</sup> The result in [4] is given for the case of multi-agent epistemic logic with a common knowledge operator, assuming  $\mathcal{M}$  is finite. The proof for the single-agent case is trivial and requires no assumption of finiteness.

and self-refuting formulas are those for which  $[\varphi]\neg\varphi$  is valid (in **S5**). While we will not deal directly with this language, our work will apply indirectly to PAL given the following result, due to Plaza [17], which relates PAL to the basic modal language.

**Theorem 2.6** *Every formula in the language of PAL is equivalent (over partitions) to a formula in the basic modal language.*

Direct results on successful formulas in PAL with multiple knowledge modalities and the *common knowledge* operator have been obtained by van Ditmarsch and Kooi in [6], which also presents an analysis of the role of (un)successful formulas in a variety of scenarios involving information change.

Another benefit of the PAL definitions of successful and self-refuting is that they give an upper bound on the complexity of checking a formula for these properties. The following result including a lower bound for the success problem is due to Johan van Benthem in correspondence.

**Theorem 2.7** *The success problem for S5 is coNP-complete.*

**Proof.** For an upper bound,  $\varphi$  is successful iff  $[\varphi]\varphi$  is valid, and the validity problem for single-agent PAL is coNP-complete [14]. For a lower bound, we use the following reduction of the validity problem for **S5**, which is coNP-complete, to the success problem for **S5**:  $\varphi$  is valid iff for a new variable  $p$ ,  $\psi \equiv p \wedge (\Diamond\neg p \vee \varphi)$  is successful. From left to right, if  $\mathcal{M}, w \models \psi$  then  $\mathcal{M}_{|\psi}, w \models p$ , and since  $\varphi$  is valid,  $\mathcal{M}_{|\psi}, w \models \varphi$ , so  $\mathcal{M}_{|\psi}, w \models \psi$ . From right to left, take any pointed model  $\mathcal{M}, w$ . Extend the language with a new variable  $p$  and extend  $\mathcal{M}$  to  $\mathcal{M}'$  with one new point  $v$ , related to all other points, with the same valuation as  $w$  for all variables of the old language. Make  $p$  true everywhere except at  $v$ . Then since  $\mathcal{M}', w \models p \wedge \Diamond\neg p$ , we have  $\mathcal{M}', w \models \psi$ , and since  $\psi$  is successful, we have  $\mathcal{M}'_{|\psi}, w \models \psi$ . But  $v$  is eliminated in  $\mathcal{M}'_{|\psi}$ , so  $\mathcal{M}'_{|\psi}, w \not\models \Diamond\neg p$ . Hence  $\mathcal{M}'_{|\psi}, w \models \varphi$ . We conclude that  $\mathcal{M}, w \models \varphi$ , for  $\mathcal{M}'_{|\psi}$  and  $\mathcal{M}$  differ only with respect to the valuation of  $p$ , which  $\varphi$  does not contain.  $\square$

A similar argument shows that  $\varphi$  is valid iff for a new variable  $p$ ,  $p \wedge (\Diamond\neg p \vee \neg\varphi)$  is self-refuting, so the self-refutation problem is coNP-complete given the upper bound from PAL. From left to right the reduction is obvious, while from right to left the same extension of  $\mathcal{M}$  to  $\mathcal{M}'$  works. Similar arguments also show that the success and self-refutation problems are PSPACE-complete for multimodal **S5**.

### 3 The Moorean Source of Unsuccessfulness

To show that Moorean phenomena are the source of all unsuccessfulness in logics for an introspective agent, we proceed via normal forms for such logics. The following normal form derives from Carnap [10], and Proposition 3.2 is standard.

**Definition 3.1** A formula is in *normal form* iff it is a disjunction of conjunctions of the form  $\delta \equiv \alpha \wedge \Box\beta_1 \wedge \dots \wedge \Box\beta_n \wedge \Diamond\gamma_1 \wedge \dots \wedge \Diamond\gamma_m$  where  $\alpha$  and each  $\gamma_i$  are conjunctions of literals and each  $\beta_i$  is a disjunction of literals.

**Proposition 3.2** *For every formula  $\varphi$ , there is a formula  $\varphi'$  in normal form such that  $\varphi$  and  $\varphi'$  are equivalent in **K45**.*

Theorem 1.7.6.4 of [15] gives the analogue of Proposition 3.2 for **S5**. Inspection of the proof shows that the necessary equivalences hold in **K45**.

We use the following notation and terminology in this section and Appendix A.

**Definition 3.3** Given  $\delta \equiv \alpha \wedge \Box\beta_1 \wedge \dots \wedge \Box\beta_n \wedge \Diamond\gamma_1 \wedge \dots \wedge \Diamond\gamma_m$  in normal form, we define  $\delta^\alpha \equiv \alpha$ ,  $\delta^{\alpha\Box} \equiv \alpha \wedge \Box\beta_1 \wedge \dots \wedge \Box\beta_n$ , and similarly for  $\delta^{\alpha\Diamond}$ ,  $\delta^\Box$ ,  $\delta^{\Box\Diamond}$ , and  $\delta^\Diamond$ .

**Definition 3.4** Where  $\chi$  is a conjunction or disjunction of literals, let  $L(\chi)$  be the set of literals in  $\chi$ . A set of literals is *open* iff no literal in the set is the negation of any of the others.

**Definition 3.5** A conjunction  $\delta \equiv \alpha \wedge \Box\beta_1 \wedge \dots \wedge \Box\beta_n \wedge \Diamond\gamma_1 \wedge \dots \wedge \Diamond\gamma_m$  in normal form is **KD45-clear** iff: (i)  $L(\alpha)$  is open; (ii) there is an open set of literals  $\{l_1, \dots, l_n\}$  with  $l_i \in L(\beta_i)$ ; and (iii) for every  $\gamma_k$  there is a set of literals  $\{l_1, \dots, l_n\}$  with  $l_i \in L(\beta_i)$  such that  $\{l_1, \dots, l_n\} \cup L(\gamma_k)$  is open. A disjunction in normal form is **KD45-clear** iff at least one of its disjuncts is **KD45-clear**.

In the following, by “clear” and “satisfiable” we mean **KD45-clear** and satisfiable in a quasi-partition, respectively. In the case of **S5-clear**, we must require in clause (ii) of Definition 3.5 that  $\{l_1, \dots, l_n\} \cup L(\alpha)$  is open, in which case (i) is unnecessary. In both cases, the following lemma holds with the appropriate definition of clarity.

**Lemma 3.6** *A formula in normal form is satisfiable iff it is clear.*

**Proof.** [Sketch] Suppose  $\varphi$  in normal form is satisfiable. Then there is some disjunct  $\delta$  of  $\varphi$  satisfied at a pointed model. Read off the appropriate open sets of literals from the current point for clarity condition (i), from some accessible point for (ii), and from witnesses for each  $\Diamond\gamma_k$  in  $\delta$  for (iii). In the other direction, suppose there are open sets of literals as described. Construct a model where  $\delta^\alpha$  is true at the root point  $w$ , which is possible by clarity condition (i). For each  $\Diamond\gamma_k$ , add a point  $v$  accessible from  $w$  and extend the valuation such that all conjuncts of  $\gamma_k$  and some disjunct of each  $\beta_i$  are true at  $v$ , which is possible by (iii). If there are no  $\Diamond\gamma_k$  formulas in  $\delta$ , add an accessible point  $v$  and extend the valuation such that some disjunct of each  $\beta_i$  is true at  $v$ , which is possible by (ii). Then  $\varphi$  is true at  $w$ .  $\square$

The following definition fixes the basic class of Moore conjunctions, as well as a wider class of Moorean conjunctions. The intuition is that a Moore conjunction simultaneously asserts a lack of information about another fact being asserted, and a Moorean conjunction is one that behaves like a Moore conjunction in some context.

We write  $\sim\varphi$  for the negation of  $\varphi$  in negation normal form, with  $\neg$  applying only to literals.

**Definition 3.7** Let  $\delta$  be a conjunction in normal form.

- (i)  $\delta$  is a *Moore conjunction* iff  $\delta \wedge \Diamond\delta^\alpha$  is not clear or there is a  $\Diamond\gamma_k$  conjunct in  $\delta$  such that  $\delta \wedge \Diamond(\delta^\alpha \wedge \gamma_k)$  is not clear.
- (ii)  $\delta$  is a *Moorean conjunction* iff there is a  $\Diamond\gamma_k$  conjunct in  $\delta$  such that  $\delta \wedge \Diamond\delta^\alpha \wedge \Box(\sim\delta^\alpha \vee \sim\gamma_k)$  is clear.

Assuming **S5**,  $\delta \wedge \Diamond\delta^\alpha$  may be replaced by  $\delta$  in both (i) and (ii).



The definition of a Moore conjunction generalizes from the paradigmatic case of  $p \wedge \Diamond \neg p$  to include formulas such as  $p \wedge \Diamond q \wedge \Box (q \rightarrow \neg p)$ . The formulas  $p \wedge \neg p$  and  $p \wedge \Box \neg p$  are also Moore conjunctions, since for these  $\delta$  the formula  $\delta \wedge \Diamond \delta^\alpha$  is not clear, but for the same reason they are not Moorean conjunctions. An example of a Moorean conjunction that is not a Moore conjunction is  $p \wedge \Diamond q$ . We consider this formula Moorean because in a context where the agent knows that  $q$  implies  $\neg p$ , so  $\Box (q \rightarrow \neg p)$  holds, learning  $p \wedge \Diamond q$  has the same effect as learning  $p \wedge \Diamond \neg p$ . By contrast, the formula  $p \wedge \Diamond q \wedge \Diamond (p \wedge q)$  rules out the Moorean context with its last conjunct, and it is not Moorean according to Definition 3.7(ii).

In the proofs of Lemma 3.9 and Theorem 3.13 below, we will assume without loss of generality that all quasi-partitions considered are *chained*, i.e.,  $\forall w, v \in W : w \neq v \Rightarrow wRv \vee vRw$ ,<sup>3</sup> in which case the following basic facts hold.

**Lemma 3.8** *Where  $\mathcal{M}$  is a chained quasi-partition and  $\delta$  is a conjunction in normal form:*

- (i)  $\mathcal{M}, w \models \delta \Rightarrow \mathcal{M} \models \delta^{\Box \Diamond}$ ;
- (ii)  $\mathcal{M}, w \models \delta \Rightarrow \mathcal{M}_{|\delta} = \mathcal{M}_{|\delta^\alpha}$  and  $\mathcal{M}_{|\delta} \models \delta^\alpha$ ;
- (iii)  $\mathcal{M}_{|\delta}, w \models \delta \Rightarrow \mathcal{M}_{|\delta} \models \delta$ .

We now prove the main lemma used in the proof of Theorem 3.13.

**Lemma 3.9** *Let  $\delta$  be a conjunction in normal form. The following hold for both KD45 and S5.*

- (i)  $\delta$  is self-refuting if and only if it is a Moore conjunction.
- (ii)  $\delta$  is unsuccessful if and only if it is a Moorean conjunction.

**Proof.** ( $\Leftarrow$  (i)) Suppose  $\delta$  is a Moore conjunction. Case 1:  $\delta \wedge \Diamond \delta^\alpha$  is not clear. By Lemmas 3.6 and 3.8(i),  $\delta \wedge \Diamond \delta^\alpha$  is clear iff  $\Diamond^+ \delta$  is satisfiable, so in this case  $\Diamond^+ \delta$  is unsatisfiable and hence  $\delta$  is self-refuting. Case 2:  $\delta \wedge \Diamond \delta^\alpha$  is clear, so suppose  $\mathcal{M}, w \models \Diamond^+ \delta$ . For the  $\Diamond \gamma_k$  in  $\delta$  such that  $\delta \wedge \Diamond (\delta^\alpha \wedge \gamma_k)$  is not clear,  $\neg \delta \vee \Box (\neg \delta^\alpha \vee \neg \gamma_k)$  is valid by Lemma 3.6, so  $\mathcal{M}, w \models \Box (\neg \delta^\alpha \vee \neg \gamma_k)$  given  $\mathcal{M}, w \models \delta$ . Then since  $\Box (\neg \delta^\alpha \vee \neg \gamma_k)$  is universal, it is preserved under submodels by Theorem 2.4, so  $\mathcal{M}_{|\delta}, w \models \Box (\neg \delta^\alpha \vee \neg \gamma_k)$ . From Lemma 3.8(ii),  $\mathcal{M}_{|\delta}, w \models \Box \delta^\alpha$ , so  $\mathcal{M}_{|\delta}, w \models \Box \neg \gamma_k$ . Since  $\Diamond \gamma_k$  is a conjunct in  $\delta$ ,  $\mathcal{M}_{|\delta}, w \not\models \delta$ . Since  $\mathcal{M}$  was arbitrary,  $\delta$  is self-refuting.

((i)  $\Rightarrow$ ) We prove the contrapositive.<sup>4</sup> Suppose  $\delta$  is not a Moore conjunction. Then  $\delta \wedge \Diamond \delta^\alpha$  is clear, and for every  $\Diamond \gamma_k$  in  $\delta$ ,  $\delta \wedge \Diamond (\delta^\alpha \wedge \gamma_k)$  is clear (\*). If there are no  $\Diamond \gamma_k$  conjuncts in  $\delta$ , then  $\delta$  is a universal formula with  $\Diamond^+ \delta$  satisfiable, so it is not self-refuting by Theorem 2.4. Suppose there are  $\Diamond \gamma_k$  conjuncts in  $\delta$ . We claim that  $\delta' \equiv \delta \wedge \Box \beta_{n+1} \wedge \dots \wedge \Box \beta_{n+j}$  is clear where  $\{\beta_{n+1}, \dots, \beta_{n+j}\} = L(\delta^\alpha)$ . Given assumption (\*) and clarity condition (iii) for each  $\delta \wedge \Diamond (\delta^\alpha \wedge \gamma_k)$ , we have that for all  $\Diamond \gamma_k$  in  $\delta$  there is a set  $\{l_1, \dots, l_n\}$  with  $l_i \in L(\beta_i)$  such that  $\{l_1, \dots, l_n\} \cup L(\delta^\alpha \wedge \gamma_k)$  is open. Taking  $\{l_{n+1}, \dots, l_{n+j}\} = L(\delta^\alpha)$ ,  $\{l_1, \dots, l_{n+j}\} \cup L(\gamma_k) = \{l_1, \dots, l_n\} \cup L(\delta^\alpha \wedge \gamma_k)$  is open, which gives clarity conditions (i), (ii) and (iii) for  $\delta'$ . Since  $\delta'$  is clear, suppose  $\mathcal{N}, w \models \delta'$ . From the fact that  $\models \delta' \leftrightarrow (\delta \wedge \Box \delta^\alpha)$  we have  $\mathcal{N}, w \models \delta \wedge \Box \delta^\alpha$ . Given

<sup>3</sup> Thanks to Dr. Yanjing Wang for catching the omission of  $w \neq v$  in the published version.

<sup>4</sup> The following argument establishes something stronger than we need for Lemma 3.9, but we use it to establish Corollary 5.3 below. Compare the ( $\Rightarrow$ ) direction of the proof of Theorem A.3 in Appendix A.

$\mathcal{N}, w \models \Box\delta^\alpha$  and the assumption that  $\mathcal{N}$  is a chained quasi-partition,  $\mathcal{N} \models \delta^\alpha$ ; given  $\mathcal{N}, w \models \delta$  and Lemma 3.8(i),  $\mathcal{N} \models \delta^{\Box\Diamond}$ . Hence  $\mathcal{N} \models \delta$ , in which case  $\mathcal{N}, w \models \Diamond^+\delta$  and  $\mathcal{N}_{|\delta} = \mathcal{N}$ . It follows that  $\mathcal{N}_{|\delta}, w \models \delta$ , so  $\delta$  is not self-refuting.

((i)  $\Leftarrow$ ) Suppose  $\delta$  is a Moorean conjunction. Since for some  $\Diamond\gamma_k$  in  $\delta$ ,  $\chi \equiv \delta \wedge \Diamond\delta^\alpha \wedge \Box(\sim\delta^\alpha \vee \sim\gamma_k)$  is clear, there is a model with  $\mathcal{M}, w \models \chi$ . Given  $\mathcal{M}, w \models \delta \wedge \Diamond\delta^\alpha$ , we have  $\mathcal{M}, w \models \Diamond^+\delta$  by Lemma 3.8(i). Given  $\mathcal{M}, w \models \Box(\sim\delta^\alpha \vee \sim\gamma_k)$ , by the same reasoning as in Case 2 of ((i)  $\Leftarrow$ ),  $\mathcal{M}_{|\delta}, w \not\models \delta$ . Therefore  $\delta$  is unsuccessful.

((i)  $\Rightarrow$ ) We prove the contrapositive. Suppose  $\mathcal{M}, w \models \Diamond^+\delta$  and  $\delta$  is not a Moorean conjunction. Then for all  $\Diamond\gamma_k$  in  $\delta$ ,  $\chi_k \equiv \delta \wedge \Diamond\delta^\alpha \wedge \Box(\sim\delta^\alpha \vee \sim\gamma_k)$  is not clear. To show  $\mathcal{M}_{|\delta}, w \models \delta$ , it suffices to show  $\mathcal{M}_{|\delta}, w \models \delta^\Diamond$ , since  $\delta^{\alpha\Box}$  is universal and therefore preserved under submodels. Consider some  $\Diamond\gamma_k$  conjunct in  $\delta$ . It follows from our assumption that  $\chi_k$  is unsatisfiable, whence  $(\delta \wedge \Diamond\delta^\alpha) \rightarrow \Diamond(\delta^\alpha \wedge \gamma_k)$  is valid. Then from  $\mathcal{M}, w \models \Diamond^+\delta$  we obtain  $\mathcal{M}, w \models \Diamond(\delta^\alpha \wedge \gamma_k)$ , so there is a  $v$  with  $wRv$  and  $\mathcal{M}, v \models (\delta^\alpha \wedge \gamma_k)$ . By Lemma 3.8(ii),  $v$  is retained in  $\mathcal{M}_{|\delta}$ . Since  $\gamma_k$  is propositional,  $\mathcal{M}_{|\delta}, v \models \gamma_k$  and hence  $\mathcal{M}_{|\delta}, w \models \Diamond\gamma_k$ . Since  $\Diamond\gamma_k$  was arbitrary,  $\mathcal{M}_{|\delta}, w \models \delta^\Diamond$  and hence  $\mathcal{M}_{|\delta}, w \models \delta$ . Since  $\mathcal{M}$  was arbitrary,  $\delta$  is successful.  $\square$

Lemma 3.9 gives necessary and sufficient conditions for the successfulness of a conjunction in normal form. We now introduce an apparently stronger notion.

**Definition 3.10** A formula  $\varphi$  is *super-successful* (in  $\mathbf{L}$ ) iff for every pointed model (of  $\mathbf{L}$ ),  $\mathcal{M}, w \models \Diamond^+\varphi$  implies  $\mathcal{M}', w \models \varphi$  for every  $\mathcal{M}'$  such that  $\mathcal{M}_{|\varphi} \subseteq \mathcal{M}' \subseteq \mathcal{M}$ .

If  $\varphi$  is super-successful and  $\mathcal{M}, w \models \varphi$ , then as points that are not in  $\mathcal{M}_{|\varphi}$  are eliminated from  $\mathcal{M}$ ,  $\varphi$  remains true at  $w$ . Since we take the elimination of points as an agent's acquisition of new information, this means that  $\varphi$  remains true as the agent approaches, by way of the incremental acquisition of new information, the epistemic state of  $\mathcal{M}_{|\varphi}$  wherein the agent knows  $\varphi$ . Intuitively, we can say that a super-successful formula remains true while an agent is “on the way” to learning it.

We will use the next lemma in the proof of Theorem 3.13.

**Lemma 3.11** *If  $\delta$  is a successful conjunction in normal form,  $\delta$  is super-successful.*

**Proof.** Suppose  $\delta$  is not super-successful, so there is a pointed model such that  $\mathcal{M}, w \models \Diamond^+\delta$  and an  $\mathcal{M}'$  such that  $\mathcal{M}_{|\delta} \subseteq \mathcal{M}' \subseteq \mathcal{M}$  and  $\mathcal{M}', w \not\models \delta$ . Since  $\mathcal{M}' \subseteq \mathcal{M}$  and  $\delta^{\alpha\Box}$  is preserved under submodels, we must have  $\mathcal{M}', w \not\models \delta^\Diamond$ . But then  $\mathcal{M}_{|\delta}, w \not\models \delta^\Diamond$  given that  $\mathcal{M}_{|\delta} \subseteq \mathcal{M}'$  and  $\delta^\Diamond$  is preserved under extensions. Hence  $\mathcal{M}_{|\delta}, w \not\models \delta$ , so  $\delta$  is unsuccessful.  $\square$

We now lift the definition of Moore and Moorean to arbitrary formulas.

**Definition 3.12** Let  $\varphi$  be an arbitrary formula.

- (i)  $\varphi$  is a *Moore sentence* iff any normal form of  $\varphi$  is a disjunction of Moore conjunctions.
- (ii)  $\varphi$  is a *Moorean sentence* iff any normal form of  $\varphi$  contains a Moorean conjunction as a disjunct.

The following theorem gives necessary conditions for self-refuting and unsuccessful formulas. In Appendix A we strengthen this result with conditions that are sufficient as well as necessary.



**Theorem 3.13** *Let  $\varphi$  be an arbitrary formula.*

- (i) *If  $\varphi$  is self-refuting in any sublogic of **S5**, then  $\varphi$  is a Moore sentence.*
- (ii) *If  $\varphi$  is unsuccessful in any extension of **KD45**, then  $\varphi$  is a Moorean sentence.*

**Proof.** By Lemma 2.3 it suffices to show the consequent of (i) for  $\varphi$  that is self-refuting in **S5** and the consequent of (ii) for  $\varphi$  that is unsuccessful in **KD45**. Since  $\varphi$  is self-refuting (resp. unsuccessful) iff any equivalent normal form of  $\varphi$  is self-refuting (resp. unsuccessful), let us assume that  $\varphi$  is already in normal form.

(i) Suppose  $\varphi$  is not a Moore sentence, so by Definition 3.12 there is a disjunct  $\delta$  of  $\varphi$  that is not a Moore conjunction. Then by Lemma 3.9,  $\delta$  is not self-refuting, so there is a pointed model with  $\mathcal{M}, w \vDash \diamond^+\delta$  and  $\mathcal{M}_{|\delta}, w \vDash \delta$ . By Lemma 3.8(iii),  $\mathcal{M}_{|\delta} \vDash \delta$ . It follows that  $(\mathcal{M}_{|\delta})_{|\varphi} = \mathcal{M}_{|\delta}$ , since all points in  $\mathcal{M}_{|\delta}$  satisfy one of the disjuncts of  $\varphi$ , namely  $\delta$ . Then given  $\mathcal{M}_{|\delta}, w \vDash \diamond^+\delta$ , we have  $\mathcal{M}_{|\delta}, w \vDash \diamond^+\varphi$  and hence  $(\mathcal{M}_{|\delta})_{|\varphi}, w \vDash \varphi$ . Therefore  $\varphi$  is not self-refuting.

(ii) We prove something stronger. Suppose  $\varphi$  is not a Moorean sentence, so by Definition 3.12 no disjunct of  $\varphi$  is a Moorean conjunction. Then each disjunct of  $\varphi$  is successful by Lemma 3.9. Consider a pointed model such that  $\mathcal{M}, w \vDash \diamond^+\varphi$ , so for some disjunct  $\delta$  of  $\varphi$ , we have  $\mathcal{M}, w \vDash \delta$ . Since  $\varphi$  is a disjunction,  $\mathcal{M}_{|\delta} \subseteq \mathcal{M}_{|\varphi}$ . By Lemma 3.11,  $\delta$  is super-successful, so for any  $\mathcal{M}'$  with  $\mathcal{M}_{|\varphi} \subseteq \mathcal{M}' \subseteq \mathcal{M}$ , we have  $\mathcal{M}', w \vDash \delta$  and hence  $\mathcal{M}', w \vDash \varphi$ . Since  $\mathcal{M}$  was arbitrary,  $\varphi$  is super-successful.  $\square$

## 4 Unsuccessfulness in Other Logics

We now consider the sources of unsuccessfulness in logics for an agent without introspection (logics without axioms 4 and 5) and in logics for multiple agents.

From an epistemic perspective, the most interesting (normal) proper sublogics of **S5** are obtained by dropping axiom 5 and adding something weaker in its place. Indeed, logics such as **S4**, **S4.x** for  $x = 2, 3, 4$ , etc., have been proposed as logics of knowledge. Call logics **L** and **L'** *comparable* if **L** is a sublogic of **L'** or *vice versa*.

**Proposition 4.1** *For any normal, proper sublogic **L** of **S5**, comparable to **S4.4**, there is a formula (consistent with **S5**) that is unsuccessful in **L** but is not Moorean.<sup>5</sup>*

**Proof.** First, we claim that  $\varphi \equiv \diamond p \wedge \diamond \neg p$  is unsuccessful in **S4.4** and hence in any sublogic of **S4.4** by Lemma 2.3. In the **S4.4** model  $\mathcal{M}$  in Figure 2,  $\varphi$  is true at the left point, but in  $\mathcal{M}_{|\varphi}$ , the right point is eliminated, so  $\varphi$  becomes false at the left point. Note that the formula is already in normal form and is not Moorean.

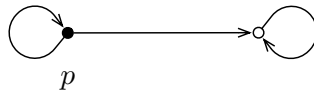


Fig. 2.

Next, we claim that  $\mathcal{M}$  is a model of any logic **L** that is a proper extension of **S4.4** and a proper sublogic of **S5**. Suppose not, so there is a theorem  $\varphi$  of **L** with  $\mathcal{M} \not\vDash \varphi$ . Change  $\varphi$  to  $\varphi'$  by substituting  $(p \wedge \neg p)$  for any propositional variable  $q$

<sup>5</sup> **S4.4** is **S4** plus  $\varphi \rightarrow (\diamond \square \varphi \rightarrow \square \varphi)$ . It is a proper extension of **S4.3** and therefore also of **S4.2** [21].

other than  $p$ . Since  $\mathbf{L}$  is normal and therefore closed under substitution,  $\varphi'$  is also a theorem of  $\mathbf{L}$  and therefore of  $\mathbf{S5}$ . Moreover, since for all variables  $q$  other than  $p$ ,  $\mathcal{M} \models \neg q$ , the substitution of  $(p \wedge \neg p)$  for  $q$  preserves (un)satisfiability in  $\mathcal{M}$ , so  $\mathcal{M} \not\models \varphi'$ . Hence  $\varphi'$  is not a theorem of  $\mathbf{S4.4}$ . But by a result of Zeman [21], for any formula  $\psi$ , containing exactly one variable, that is a theorem of  $\mathbf{S5}$  but not a theorem of  $\mathbf{S4.4}$ , adding  $\psi$  to  $\mathbf{S4.4}$  gives  $\mathbf{S5}$ . Hence  $\mathbf{L}$  is  $\mathbf{S5}$ , a contradiction. Since  $\mathcal{M}$  models any logic between  $\mathbf{S4.4}$  and  $\mathbf{S5}$ ,  $\varphi$  is unsuccessful in these logics.  $\square$

Proposition 4.1 shows that  $\mathbf{S5}$  is unique among the typical logics of knowledge insofar as all of its unsuccessful formulas are Moorean. The counterexample for the weaker logics shows that without negative introspection, one can come to know  $p$  by being truly told, “You do not know whether or not  $p$ ,” a surprising case of unsuccessfulness. The following proposition, although weaker than Proposition 4.1, shows that non-Moorean unsuccessful formulas appear if we weaken logics of knowledge and belief in other ways as well.

**Proposition 4.2** *For any sublogic  $\mathbf{L}$  of  $\mathbf{KTB}$  or  $\mathbf{KD5}$ , there is a formula (consistent with  $\mathbf{S5}$ ) that is unsuccessful in  $\mathbf{L}$  but is not Moorean.*

**Proof.** For  $\mathbf{KTB}$  consider  $\varphi \equiv \diamond p \wedge \diamond q$  and the model  $\mathcal{M}$  in Figure 3. In  $\mathcal{M}_{|\varphi}$ , the left and right points are eliminated, so  $\varphi$  becomes false at the center point.

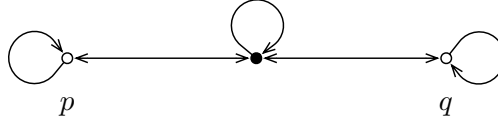


Fig. 3.

For  $\mathbf{KD5}$  consider  $\psi \equiv \neg q \vee (\Box p \wedge \diamond \diamond q)$  and the model  $\mathcal{M}$  in Figure 4. Only the right point is eliminated in  $\mathcal{M}_{|\psi}$ , so  $\psi$  becomes false at the left point.

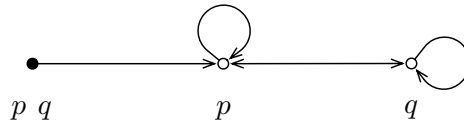


Fig. 4.

The formula  $\psi$  is equivalent in  $\mathbf{KD45}$  to the normal form  $\neg q \vee (\Box p \wedge \diamond q)$ , which is not Moorean, since neither disjunct is a Moorean conjunction.<sup>6</sup>  $\square$

There are other formulas that one may wish to categorize as Moorean, perhaps even as Moore sentences, but which are inconsistent with  $\mathbf{KD45}$ .

**Example 4.3** The formula  $\Box p \wedge \diamond \neg \Box p$  is a kind of “higher order” Moore sentence, which is satisfiable on intransitive frames but is also self-refuting over such frames. Similarly, the formula  $\diamond p \wedge \diamond \neg \diamond p$  is satisfiable yet self-refuting on non-Euclidean frames. The first formula says that the agent is not aware of what he believes, while the second says that he is not aware of what he does not believe. These formulas are the very witnesses of a failure to validate axioms 4 and 5, respectively.

<sup>6</sup> Note that while the right disjunct is not a Moorean conjunction, in  $\mathbf{S5}$  (but not in  $\mathbf{KD45}$ ) it implies the Moorean conjunction  $p \wedge \diamond q$ . This is an instance of the more general fact that in  $\mathbf{S5}$  some successful formulas imply unsuccessful formulas.

A natural question is whether there are non-Moorean sources of unsuccessfulness in languages more expressive than the basic modal language. Consider, for example, a language with multiple modalities, as in multi-agent epistemic logic. Without giving a formal definition of Moorean in the multi-agent case, it is nonetheless clear that there are more ways to be Moorean in the multi-agent case than in the single-agent case, even assuming quasi-partitions for each relation. For example, there are self-refuting, *indirect* Moore sentences such as  $\Box_a p \wedge \Diamond_b \neg p$ , which imply single-agent Moore sentences (in this case,  $p \wedge \Diamond_b \neg p$ ) in multi-agent **S5**. There are also self-refuting, *higher order* Moore sentences such as  $\Box_a p \wedge \Diamond_b \neg \Box_a p$  and  $\Diamond_a p \wedge \Diamond_b \neg \Diamond_a p$ , which resemble the higher order, single-agent Moore sentences consistent with logics weaker than **KD45**, noted in Example 4.3. However, not all unsuccessful formulas in the multi-agent context have a Moorean character.

**Example 4.4** In the model  $\mathcal{M}$  in Figure 5, the formula  $\varphi \equiv \Diamond_a p \wedge \Diamond_a \Diamond_b \neg p$  is true at the left point but false at the right point, since  $\Diamond_a p$  is false at the right point.

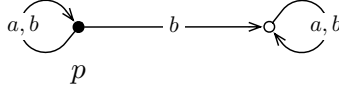


Fig. 5.

As a result, in the model  $\mathcal{M}_{|\varphi}$  the right point is eliminated, in which case  $\varphi$  becomes false at the left point because  $\Diamond_a \Diamond_b \neg p$  becomes false there.

The formula  $\Diamond_a p \wedge \Diamond_a \Diamond_b \neg p$  does not resemble any single-agent Moorean formula, yet it is nonetheless unsuccessful. Given the connection that we have observed between introspection and Moorean phenomena in the single-agent case, we can see why there should be non-Moorean unsuccessful formulas in the multi-agent case: agents do not have introspective access to each other's knowledge or beliefs.

## 5 Related Classes of Formulas

Theorem 3.13 shows that in a wide range of logics, every self-refuting formula is a Moore sentence (i) and every unsuccessful formula is a Moorean sentence (ii). However, neither the converse of (i) nor the converse of (ii) holds in general. In this section, we relate the failures of the converses of (i) and (ii) to other classes of formulas. In Theorems A.3 and A.6 in Appendix A, we overcome these failures and give syntactic characterizations of the self-refuting and unsuccessful formulas as the *strong* Moore sentences and *strong* Moorean sentences, respectively.

For simplicity we assume in this section that we are working with **S5** only.

### 5.1 Informative, Cartesian, and eventually self-refuting formulas

**Definition 5.1** A formula  $\varphi$  is (*potentially*) *informative* iff there is a pointed model such that  $\mathcal{M}, w \models \varphi$  and  $\mathcal{M}_{|\varphi} \neq \mathcal{M}$ . Otherwise  $\varphi$  is *uninformative*. A formula  $\varphi$  is *always informative* iff for all pointed models such that  $\mathcal{M}, w \models \varphi$ ,  $\mathcal{M}_{|\varphi} \neq \mathcal{M}$ .

Note that if a formula is not always informative, then it is not self-refuting, for there is a model such that  $\mathcal{M}, w \models \varphi$  but  $\mathcal{M}_{|\varphi} = \mathcal{M}$ , so  $\mathcal{M}_{|\varphi}, w \models \varphi$ . This observation

explains some of the counterexamples to the converse of Theorem 3.13(i), Moore sentences that are not self-refuting.

**Example 5.2** The formula  $(p \wedge \Diamond \neg p) \vee (q \wedge \Diamond \neg q)$  is a Moore sentence, but it is not always informative (take a model with only two connected points, one of which satisfies  $p$  but not  $q$  and the other of which satisfies  $q$  but not  $p$ ) and hence not self-refuting. The formula  $(p \wedge \Diamond \neg p) \vee (\neg p \wedge \Diamond p)$  is a Moore sentence, but it is uninformative and hence successful, but not self-refuting. These examples show that neither the self-refuting nor the unsuccessful formulas are closed under disjunction.

From our previous results, we have the following corollary.

**Corollary 5.3** *A conjunction in normal form is always informative iff it is self-refuting.*

**Proof.** By inspection of the proof of Lemma 3.9(i).  $\square$

However, a Moore sentence that is always informative may not be self-refuting.

**Example 5.4** The formula  $\varphi \equiv (p \wedge \Diamond \neg p) \vee (p \wedge q \wedge \Diamond \neg q)$  is always informative, for if  $\mathcal{M}, w \models (p \wedge \Diamond \neg p)$ , then there is a witness to  $\Diamond \neg p$  that is eliminated in  $\mathcal{M}_{|\varphi}$ , and if  $\mathcal{M}, w \models (p \wedge q \wedge \Diamond \neg q)$ , then either the witness to  $\Diamond \neg q$  does not satisfy  $(p \wedge \Diamond \neg p)$ , in which case it does not satisfy  $\varphi$  and is eliminated in  $\mathcal{M}_{|\varphi}$ , or it does satisfy  $(p \wedge \Diamond \neg p)$ , in which case there is a witness to  $\Diamond \neg p$  that is eliminated in  $\mathcal{M}_{|\varphi}$ . In either case,  $\mathcal{M}_{|\varphi} \neq \mathcal{M}$ , so  $\varphi$  is always informative. However,  $\varphi$  is not self-refuting, as shown by the partition model  $\mathcal{M}$  with  $W = \{w, v, u\}$ ,  $V(p) = \{w, v\}$ , and  $V(q) = \{w, u\}$ . We have  $\mathcal{M}, w \models p \wedge q \wedge \Diamond \neg q$ , and given  $W_{|\varphi} = \{w, v\}$ , also  $\mathcal{M}_{|\varphi}, w \models p \wedge q \wedge \Diamond \neg q$ .

Interestingly, the formula  $\varphi$  is *self-refuting within two steps*; given  $\mathcal{M}, w \models \varphi$  and  $\mathcal{M}_{|\varphi}, w \models \varphi$ , we have  $(\mathcal{M}_{|\varphi})_{|\varphi}, w \not\models \varphi$ . Let  $\delta_1 \equiv (p \wedge \Diamond \neg p)$  and  $\delta_2 \equiv (p \wedge q \wedge \Diamond \neg q)$  be the disjuncts of  $\varphi$ . All  $\neg p$ -points are eliminated in  $\mathcal{M}_{|\varphi}$ , so  $\mathcal{M}_{|\varphi} \models \neg \delta_1$ . Hence  $(\mathcal{M}_{|\varphi})_{|\varphi} = (\mathcal{M}_{|\varphi})_{|\delta_2}$ . But since  $\mathcal{M}_{|\varphi}, w \models \delta_2$  and  $\delta_2$  is self-refuting,  $(\mathcal{M}_{|\varphi})_{|\delta_2}, w \not\models \delta_2$ . Since  $\delta_1$  is existential,  $(\mathcal{M}_{|\varphi})_{|\delta_2}, w \not\models \delta_1$ , so we conclude that  $(\mathcal{M}_{|\varphi})_{|\varphi}, w \not\models \varphi$ .

Example 5.4 points to the interest of self-refutation “in the long run.”

**Definition 5.5** Given a model  $\mathcal{M}$ , we define  $\mathcal{M}_{|n\varphi}$  recursively by  $\mathcal{M}_{|0\varphi} = \mathcal{M}$ ,  $\mathcal{M}_{|n+1\varphi} = (\mathcal{M}_{|n\varphi})_{|\varphi}$ . A formula  $\varphi$  is *self-refuting within  $n$  steps* iff for all pointed models, if  $\mathcal{M}, w \models \varphi$ , then  $\mathcal{M}_{|m\varphi}, w \not\models \varphi$  for some  $m \leq n$ ;  $\varphi$  is *eventually self-refuting* iff for all pointed models, if  $\mathcal{M}, w \models \varphi$ , then there is an  $n$  such that  $\mathcal{M}_{|n\varphi}, w \not\models \varphi$ .

Using Definition 5.5, we can generalize Corollary 5.3. First, recall Fitch’s paradox, mentioned in Section 1, which we can restate as: *if all truths are knowable, then all truths are (already) known*. One proposal for avoiding the paradox [19] is to restrict the claim that all truths are knowable to the claim that all *Cartesian* truths are knowable, in which case it does not follow that all truths are known.

**Definition 5.6**  $\varphi$  is *Cartesian* iff  $\Box\varphi$  is satisfiable.<sup>7</sup>

<sup>7</sup> The term ‘Cartesian’ is due to Tennant [19], though his definition is in terms of consistency rather than satisfiability. The term ‘knowable’ seems more natural, but ‘knowable’ is used in a different sense in the literature [9], for what we call ‘learnable’ in Definition 5.13.

The class of formulas  $\varphi$  for which  $\Box\varphi$  is satisfiable seems a natural object of study in its own right. The following proposition establishes a connection between these formulas, defined “statically,” and the other formulas classes that we have defined “dynamically” in terms of model transformations. It also generalizes Corollary 5.3.

**Proposition 5.7** *The following are equivalent:*

- (i)  $\varphi$  is always informative.
- (ii)  $\varphi$  is not Cartesian.
- (iii)  $\varphi$  is eventually self-refuting.

**Proof.** Suppose  $\varphi$  is always informative and  $\mathcal{M}, w \models \varphi$ . Without loss of generality, assume  $\mathcal{M}$  is chained (as above Lemma 3.8). Since  $\varphi$  is always informative,  $\mathcal{M}_{|\varphi} \neq \mathcal{M}$ . Hence there is a point  $v$  such that  $wRv$  and  $\mathcal{M}, v \not\models \varphi$ , in which case  $\mathcal{M}, w \not\models \Box\varphi$ . Since  $\mathcal{M}$  was arbitrary,  $\Box\varphi$  is unsatisfiable, so  $\varphi$  is not Cartesian.

Suppose  $\varphi$  is not Cartesian, and without loss of generality, assume  $\varphi$  is in normal form. We prove by induction that for all  $n \geq 0$ , if  $\mathcal{M}_{|n\varphi}, w \models \varphi$ , then there are  $\Diamond\gamma_1, \dots, \Diamond\gamma_{n+1}$  distinct subformulas of  $\varphi$  with  $\mathcal{M}_{|n+1\varphi} \models \neg\Diamond\gamma_1 \wedge \dots \wedge \neg\Diamond\gamma_{n+1}$  (\*). It follows that since  $\varphi$  has some finite number  $n$  of distinct  $\Diamond\gamma_k$  subformulas, we must have  $\mathcal{M}_{|m\varphi}, w \not\models \varphi$  for some  $m \leq n$ . For the base case, assume  $\mathcal{M}, w \models \varphi$ . Since  $\varphi$  is not Cartesian,  $\mathcal{M}_{|\varphi}, w \not\models \Box\varphi$ . Let  $v$  be such that  $wRv$  and  $\mathcal{M}_{|\varphi}, v \not\models \varphi$ , and note that  $\mathcal{M}, v \models \varphi$ , for otherwise  $v$  would not have been retained in  $\mathcal{M}_{|\varphi}$ . From the fact that universal formulas are preserved under submodels, it follows that for some  $\Diamond\gamma_1$  in  $\varphi$ , we have  $\mathcal{M}, v \models \Diamond\gamma_1$  but  $\mathcal{M}_{|\varphi}, v \not\models \Diamond\gamma_1$ . Hence  $\mathcal{M}_{|\varphi} \models \neg\Diamond\gamma_1$  by Lemma 3.8(i). For the inductive step, assume that  $\mathcal{M}_{|n+1\varphi}, w \models \varphi$ . Since  $w$  is retained in  $\mathcal{M}_{|n+1\varphi}$ , we have  $\mathcal{M}_{|n\varphi}, w \models \varphi$ , so by the induction hypothesis there are  $\Diamond\gamma_1, \dots, \Diamond\gamma_{n+1}$  distinct subformulas of  $\varphi$  for which (\*) holds. By the same reasoning as before, since  $\varphi$  is Cartesian,  $\mathcal{M}_{|n+2\varphi}, w \not\models \Box\varphi$ , so there is some  $z$  with  $wRz$  and some  $\Diamond\gamma_{n+2}$  in  $\varphi$  such that  $\mathcal{M}_{|n+1\varphi}, z \models \Diamond\gamma_{n+2}$  but  $\mathcal{M}_{|n+2\varphi} \models \neg\Diamond\gamma_{n+2}$ . Moreover, since  $\Diamond\gamma_k$  formulas are preserved under extensions,  $\mathcal{M}_{|n+2\varphi} \models \neg\Diamond\gamma_1 \wedge \dots \wedge \neg\Diamond\gamma_{n+1}$  as well. Finally, given  $\mathcal{M}_{|n+1\varphi}, z \models \Diamond\gamma_{n+2}$  and (\*),  $\Diamond\gamma_{n+2}$  is distinct from  $\Diamond\gamma_1, \dots, \Diamond\gamma_{n+1}$ .

Suppose  $\varphi$  is not always informative. Then there is a model with  $\mathcal{M}, w \models \varphi$  and  $\mathcal{M}_{|\varphi} = \mathcal{M}$ , in which case  $\mathcal{M}_{|n\varphi} = \mathcal{M}$  for all  $n$ , so  $\varphi$  is not eventually self-refuting.  $\square$

**Corollary 5.8**  *$\varphi$  is eventually self-refuting iff it is self-refuting within  $n$  steps, with  $n$  bounded by the number of distinct diamond formulas in a normal form of  $\varphi$ .*

**Proof.** By inspection of the proof of the previous proposition.  $\square$

## 5.2 Successful, super-successful, and learnable formulas

The converse of Theorem 3.13(ii) fails because a formula  $\varphi$  in normal form that contains a Moorean conjunction as a disjunct may be successful. For example, it is easy to see that if the disjunction of the non-Moorean conjunctions in  $\varphi$  is a consequence of the disjunction of the Moorean conjunctions in  $\varphi$ , then  $\varphi$  is not only successful but super-successful, given Lemma 3.11. Examples of such  $\varphi$  include  $(p \wedge \Diamond\neg p) \vee p$  and  $(p \wedge \Diamond\neg p) \vee \Diamond\neg p$ . There are also successful formulas  $\psi$  that contain Moorean conjunctions as disjuncts, but do not meet the condition of  $\varphi$ . These formulas nevertheless manage to be successful by a kind of compensation: when one disjunct

of  $\psi$  goes from true at  $\mathcal{M}, w$  to false at  $\mathcal{M}_{|\psi}, w$ , another disjunct compensates by going from false at  $\mathcal{M}, w$  to true at  $\mathcal{M}_{|\psi}, w$ . The formula  $(p \wedge \diamond \neg p) \vee \Box p$  exhibits this kind of compensation against a Moore conjunction, while the formula in the proof of the following proposition does so against a Moorean conjunction.

**Proposition 5.9** *Not all successful formulas are super-successful.*

**Proof.** The formula  $\varphi \equiv (p \wedge \diamond q) \vee \Box p$  is successful but not super-successful. Suppose  $\mathcal{M}, w \vDash \varphi$ , and without loss of generality assume that  $\mathcal{M}$  is chained. Case 1:  $\mathcal{M}, w \vDash \Box p$ . Then  $\mathcal{M}_{|\varphi} = \mathcal{M}$ , so  $\mathcal{M}_{|\varphi}, w \vDash \varphi$ . Case 2:  $\mathcal{M}, w \not\vDash \Box p$ . Then  $\mathcal{M}_{|\varphi} = \mathcal{M}_{|p \wedge \diamond q}$ , and given  $\mathcal{M}_{|p \wedge \diamond q} \vDash \Box p$ , we again have  $\mathcal{M}_{|\varphi}, w \vDash \varphi$ . Hence  $\varphi$  is successful. But  $\varphi$  is not super-successful, as shown by the partition model  $\mathcal{N}$  with  $W = \{w, v, u\}$ ,  $V(p) = \{w\}$ , and  $V(q) = \{u\}$ . We begin with  $\mathcal{N}, w \vDash \varphi$ , for while  $\mathcal{N}, w \not\vDash \Box p$ , it holds that  $\mathcal{N}, w \vDash p \wedge \diamond q$ . Moreover, given  $W_{|\varphi} = \{w\}$ , we still have  $\mathcal{N}_{|\varphi}, w \vDash \varphi$ , for while  $\mathcal{N}_{|\varphi}, w \not\vDash p \wedge \diamond q$ , it holds that  $\mathcal{N}_{|\varphi}, w \vDash \Box p$ . However, in the extension  $\mathcal{N}'$  of  $\mathcal{N}_{|\varphi}$  with  $W' = \{w, v\}$ , we now have  $\mathcal{N}', w \not\vDash \varphi$ . The fact that  $u \notin W'$  gives  $\mathcal{N}', w \not\vDash p \wedge \diamond q$ , and the fact that  $v \in W'$  gives  $\mathcal{N}', w \not\vDash \Box p$ .  $\square$

**Proposition 5.10** *A formula  $\varphi$  is super-successful iff for a propositional variable  $p$  that does not occur in  $\varphi$ ,  $\varphi \vee p$  is successful.*

**Proof.** ( $\Rightarrow$ ) Suppose  $\varphi$  is super-successful and  $\mathcal{M}, w \vDash \varphi \vee p$ , where  $p$  does not occur in  $\varphi$ . If  $\mathcal{M}, w \vDash p$ , then  $\mathcal{M}_{|\varphi \vee p}, w \vDash p$ , and if  $\mathcal{M}, w \vDash \varphi$ , then given  $\mathcal{M}_{|\varphi} \subseteq \mathcal{M}_{|\varphi \vee p}$  and the assumption that  $\varphi$  is super-successful,  $\mathcal{M}_{|\varphi \vee p}, w \vDash \varphi$ . In either case,  $\mathcal{M}_{|\varphi \vee p}, w \vDash \varphi \vee p$ , so  $\varphi \vee p$  is successful.

( $\Leftarrow$ ) Suppose  $\varphi$  is not super-successful, so there is an  $\mathcal{M} = \langle W, R, V \rangle$  with  $w \in W$  such that  $\mathcal{M}, w \vDash \varphi$ , and an  $\mathcal{M}' = \langle W', R', V' \rangle$  such that  $\mathcal{M}_{|\varphi} \subseteq \mathcal{M}' \subseteq \mathcal{M}$  and  $\mathcal{M}', w \not\vDash \varphi$ . Let  $\mathcal{M}^* = \langle W, R, V^* \rangle$  be the same as  $\mathcal{M}$  except that  $V^*(p) = W' \setminus W_{|\varphi}$ . Then since  $p$  does not occur in  $\varphi$ , we have  $\mathcal{M}^*, w \vDash \varphi \vee p$  given  $\mathcal{M}, w \vDash \varphi$ . Moreover, since  $W_{|\varphi \vee p}^* = W'$ , we have  $\mathcal{M}_{|\varphi \vee p}^*, w \not\vDash \varphi$  given  $\mathcal{M}', w \not\vDash \varphi$ . From the assumption that  $\mathcal{M}, w \vDash \varphi$ , it follows that  $w \in W_{|\varphi}$ , in which case  $w \notin V^*(p)$  and hence  $\mathcal{M}_{|\varphi \vee p}^*, w \not\vDash p$ . But then  $\mathcal{M}_{|\varphi \vee p}^*, w \not\vDash \varphi \vee p$ , so  $\varphi \vee p$  is unsuccessful.  $\square$

By Proposition 5.10, complexity and syntactic characterization results for successful formulas carry over immediately to super-successful formulas.

From the previous propositions we obtain a surprising failure of closure.

**Corollary 5.11** *The set of successful formulas is not closed under disjunction.*

**Proof.** Immediate from Propositions 5.9 and 5.10.  $\square$

**Example 5.12** The formula  $\neg p \wedge \neg q$  is successful, and by the proof of Proposition 5.9,  $(p \wedge \diamond q) \vee \Box p$  is also successful. However, the disjunction of these formulas,  $\chi \equiv (\neg p \wedge \neg q) \vee ((p \wedge \diamond q) \vee \Box p)$  is unsuccessful, as shown by the model  $\mathcal{N}$  in the proof of Proposition 5.9. We begin with  $\mathcal{N}, w \vDash \chi$ , since while  $\mathcal{N}, w \not\vDash \neg p \wedge \neg q$  (and hence  $\mathcal{N}_{|\chi}, w \not\vDash \neg p \wedge \neg q$ ), we have already seen that  $\mathcal{N}, w \vDash (p \wedge \diamond q) \vee \Box p$ . However,  $\mathcal{N}_{|\chi} = \mathcal{N}'$ , and we have already seen that  $\mathcal{N}', w \not\vDash (p \wedge \diamond q) \vee \Box p$ .

By contrast, the set of super-successful formulas is closed under disjunction, since  $\mathcal{M}_{|\varphi \vee \psi}$  is an extension of  $\mathcal{M}_{|\varphi}$  and  $\mathcal{M}_{|\psi}$ .



Another consequence of the previous results concerns the relation between successful and *learnable* formulas. Like the Cartesian formulas, the learnable formulas have been discussed in connection with Fitch’s paradox [3,9].

**Definition 5.13** A formula  $\varphi$  is (*always*) *learnable* iff for all pointed models, if  $\mathcal{M}, w \models \varphi$ , then there is some  $\psi$  such that  $\mathcal{M}_{|\psi}, w \models \Box\varphi$ .<sup>8</sup>

Formulas that are learnable (and satisfiable) are Cartesian according to Definition 5.6, but the converse does not hold [3]. For example, the formula  $p \wedge \Diamond q$  is Cartesian, since  $\Box(p \wedge \Diamond q)$  is satisfiable, but it is not (always) learnable, for if  $\mathcal{M}, w \models \Box(p \rightarrow \neg q)$  and  $\mathcal{M}' \subseteq \mathcal{M}$ , then  $\mathcal{M}', w \not\models \Box(p \wedge \Diamond q)$ .

All successful formulas are learnable [9], since if  $\mathcal{M}, w \models \varphi$  and  $\varphi$  is successful, then not only  $\mathcal{M}_{|\varphi}, w \models \varphi$  but also  $\mathcal{M}_{|\varphi}, w \models \Box\varphi$ . For if  $v$  is retained in  $\mathcal{M}_{|\varphi}$ , then  $\mathcal{M}, v \models \varphi$ , in which case  $\mathcal{M}_{|\varphi}, v \models \varphi$  by the successfulness of  $\varphi$ . Hence for a successful  $\varphi$ , we can always take  $\psi$  in Definition 5.13 to be  $\varphi$  itself. On the other hand, it is natural to ask whether some unsuccessful formulas are learnable as well.

**Corollary 5.14** *Not all learnable formulas are successful.*

**Proof.** Let  $\delta_1 \vee \delta_2$  be an unsuccessful disjunction with  $\delta_1$  and  $\delta_2$  successful, as given by Corollary 5.11. If  $\mathcal{M}, w \models \delta_1 \vee \delta_2$ , then  $\mathcal{M}, w \models \delta_i$  for  $i = 1$  or  $i = 2$ . Since  $\delta_i$  is successful, we have  $\mathcal{M}_{|\delta_i}, w \models \Box\delta_i$ , in which case  $\mathcal{M}_{|\delta_i}, w \models \Box(\delta_1 \vee \delta_2)$ . Since  $\mathcal{M}$  was arbitrary,  $\delta_1 \vee \delta_2$  is (always) learnable.  $\square$

## 6 Conclusion

From a technical point of view, we have studied the question of when a modal formula is preserved under relativizing models to the formula itself. In the epistemic interpretation of modal logic, this preservation question takes on a new significance: it concerns whether an agent retains knowledge of what is learned, which requires that what is learned remain true. We have shown (Theorem 3.13) that for an introspective agent, the only true sentences that may become false when learned are variants of the Moore sentence. For an agent without introspection or multiple agents without introspective access to each other’s knowledge or beliefs, there are non-Moorean sources of unsuccessfulness (Propositions 4.1 and 4.2, Example 4.4).

In connection with our study of Moorean phenomena, we have observed a number of related results. We saw that the sentences that always provide information to an agent, no matter the agent’s prior epistemic state, are exactly those sentences that cannot be known—and will eventually become false if repeated enough (Proposition 5.7); we saw that there are sentences that always remain true when they are learned, but whose truth value may oscillate while an agent is on the way to learning them (Proposition 5.9); and we saw that there are sentences that sometimes become false when learned directly, but which an agent can always come to know indirectly by learning something else (Corollary 5.14).

In Appendix A, we return to the problem with which we began. We give syntactic characterizations of the self-refuting and unsuccessful formulas as the *strong* Moore sentences and *strong* Moorean sentences, respectively.

<sup>8</sup> The term ‘learnable’ is used by van Benthem [3]. Balbiani et al. [9] use the term ‘knowable’.

## Acknowledgement

We wish to thank Hans van Ditmarsch and the anonymous referees for their helpful comments on an earlier draft of this paper, and Johan van Benthem for stimulating our interest in the topic of successful formulas.

## References

- [1] Baltag, A., L. Moss and S. Solecki, *The logic of public announcements, common knowledge and private suspicions*, in: I. Gilboa, editor, *Proceedings of the 7th Conference on Theoretical Aspects of Rationality and Knowledge (TARK 98)*, Morgan Kaufmann, 1998 pp. 43–56.
- [2] Baltag, A., H. van Ditmarsch and L. Moss, *Epistemic logic and information update*, in: P. Adriaans and J. van Benthem, editors, *Philosophy of Information*, North-Holland, 2008 pp. 361–456.
- [3] van Benthem, J., *What one may come to know*, *Analysis* **64** (2004), pp. 95–105.
- [4] van Benthem, J., *One is a lonely number: on the logic of communication*, in: Z. Chatzidakis, P. Koepke and W. Pohlers, editors, *Logic Colloquium '02*, ASL & A.K. Peters, 2006 pp. 96–129.
- [5] van Benthem, J., *Open problems in logical dynamics*, in: D. Gabbay, S. Goncharov and M. Zakharyashev, editors, *Mathematical Problems from Applied Logic I*, Springer, 2006 pp. 137–192.
- [6] van Ditmarsch, H. and B. Kooi, *The secret of my success*, *Synthese* **151** (2006), pp. 201–232.
- [7] van Ditmarsch, H., W. van der Hoek and B. Kooi, “Dynamic Epistemic Logic,” Springer, 2008.
- [8] Andréka, H., I. Németi and J. van Benthem, *Modal languages and bounded fragments of predicate logic*, *Journal of Philosophical Logic* **27** (1998), pp. 217–274.
- [9] Balbiani, P., A. Baltag, H. van Ditmarsch, A. Herzig, T. Hoshi and T. de Lima, ‘*Knowable*’ as ‘*known after an announcement*’, *The Review of Symbolic Logic* **1** (2008), pp. 305–334.
- [10] Carnap, R., *Modalities and quantification*, *The Journal of Symbolic Logic* **11** (1946), pp. 33–64.
- [11] Fitch, F. B., *A logical analysis of some value concepts*, *The Journal of Symbolic Logic* **28** (1963), pp. 135–142.
- [12] Gerbrandy, J., “Bisimulations on Planet Kripke,” Ph.D. thesis, University of Amsterdam (1999), ILLC Dissertation Series DS-1999-01.
- [13] Hintikka, J., “Knowledge and Belief: An Introduction to the Logic of the Two Notions,” College Publications, 2005.
- [14] Lutz, C., *Complexity and succinctness of public announcement logic*, in: P. Stone and G. Weiss, editors, *Proceedings of the Fifth International Joint Conference on Autonomous Agents and Multiagent Systems (AAMAS’06)* (2006), pp. 137–143.
- [15] Meyer, J.-J. Ch. and W. van der Hoek, “Epistemic Logic for AI and Computer Science,” *Cambridge Tracts in Theoretical Computer Science* **41**, Cambridge University Press, 1995.
- [16] Moore, G. E., *A reply to my critics*, in: P. Schilpp, editor, *The Philosophy of G.E. Moore*, The Library of Living Philosophers **4**, Northwestern University, 1942 pp. 535–677.
- [17] Plaza, J., *Logics of public communications*, in: M. Emrich, M. Pfeifer, M. Hadzikadic and Z. Ras, editors, *Proceedings of the 4th International Symposium on Methodologies for Intelligent Systems: Poster Session Program*, Oak Ridge National Laboratory, ORNL/DSRD-24, 1989 pp. 201–216.
- [18] Qian, L., *Sentences true after being announced*, in: *Proceedings of the Student Session of the 1st North American Summer School in Logic, Language, and Information (NASSLLI)*, Stanford University, 2002.
- [19] Tennant, N., “The Taming of The True,” Oxford: Clarendon Press, 1997.
- [20] Visser, A., J. van Benthem, D. de Jongh and G. R. R. de Lavalette, *NNIL, A study in intuitionistic propositional logic*, Logic Group Preprint Series 111, Philosophical Institute, Utrecht University (1994).
- [21] Zeman, J. J., *A study of some systems in the neighborhood of S4.4*, *Notre Dame Journal of Formal Logic* **12** (1971), pp. 341–357.

## A Appendix

In this appendix, we address the problem of giving syntactic characterizations of the successful and self-refuting formulas. Given Theorem 2.6, Lemma 3.6, and the PAL definitions of successful and self-refuting, there is a trivial characterization of both classes:  $\varphi$  is successful iff the result of reducing the PAL formula  $\neg[\varphi]\varphi$  to normal form in the basic modal language is not clear;  $\varphi$  is self-refuting iff the result of reducing the PAL formula  $\neg[\varphi]\neg\varphi$  to normal form in the basic modal language is not clear. Although by Proposition 2.7, the PAL definitions of successful and self-refuting lead to optimal methods for checking for these properties, they do not provide much insight beyond the semantic definitions of the formula classes. In Theorems A.3 and A.6 below, we give syntactic characterizations of the self-refuting and unsuccessful formulas that reveal more than the trivial characterization.

We state the following results for **S5** with the standard definitions of successful ( $\forall \mathcal{M}, w : \mathcal{M}, w \models \varphi \Rightarrow \mathcal{M}|_\varphi, w \models \varphi$ ) and self-refuting ( $\forall \mathcal{M}, w : \mathcal{M}, w \models \varphi \Rightarrow \mathcal{M}|_\varphi, w \not\models \varphi$ ). With minor changes the results also hold for **KD45**, given the modified definitions using precondition  $\diamond^+\varphi$  instead of  $\varphi$ .

We will continue to use the abbreviations  $\delta^\alpha$ ,  $\delta^{\square\diamond}$ , etc., from Definition 3.3. We will also use  $\hat{\delta}^{\square\diamond}$  to denote a conjunct of  $\delta^{\square\diamond}$ , i.e., a  $\square\beta_i$  or  $\diamond\gamma_k$  in  $\delta$ . As before,  $\sim\varphi$  is the negation of  $\varphi$  in negation normal form, with  $\neg$  applying only to literals.

**Definition A.1**  $\varphi$  is a *strong Moore sentence* iff for any normal form  $\varphi^*$  of  $\varphi$ , no disjunct of  $\varphi^*$  is *compensated in*  $\varphi^*$ . A disjunct  $\delta \equiv \alpha \wedge \square\beta_1 \wedge \dots \wedge \square\beta_n \wedge \diamond\gamma_1 \wedge \dots \wedge \diamond\gamma_m$  of  $\varphi^*$  is compensated in  $\varphi^*$  iff there is a disjunct  $\delta'$  of  $\varphi^*$ , a subset  $S$  of the disjuncts of  $\varphi^*$ , a sequence of disjuncts  $\sigma_{\gamma_1}, \dots, \sigma_{\gamma_m}$  (not necessarily distinct) from  $S$ , and for every  $\sigma \notin S$  a conjunct  $\hat{\sigma}^{\square\diamond}$  of  $\sigma^{\square\diamond}$ , such that  $\chi \equiv \delta' \wedge \delta^\alpha \wedge \chi_\square \wedge \chi_\diamond$  is clear, where

$$\chi_\square \equiv \bigwedge_{\sigma \in S} \sigma^{\square\diamond} \wedge \bigwedge_{\sigma \notin S} \sim \hat{\sigma}^{\square\diamond} \wedge \bigwedge_{\sigma \in S, j \leq n} \square(\sim \sigma^\alpha \vee \beta_j);$$

$$\chi_\diamond \equiv \bigwedge_{k \leq m} \diamond(\sigma_{\gamma_k}^\alpha \wedge \gamma_k).$$

Note that  $\chi$  is in normal form, so the definition of clarity properly applies. The reason for the use of  $\sim \hat{\sigma}^{\square\diamond}$  is that  $\sim \sigma^{\square\diamond}$  may be a disjunction, in which case  $\chi$  would not be in normal form. Hence we pull the disjunction out of the formula and add an existential quantifier to the condition, since  $\sigma^{\square\diamond}$  is false iff there is some conjunct  $\hat{\sigma}^{\square\diamond}$  of  $\sigma^{\square\diamond}$  that is false.

As we will see in the proof of Theorem A.3, the existence of a compensated disjunct in  $\varphi$  is equivalent to there being a model in which  $\varphi$  has a “successful update,” in the sense that  $\mathcal{M}, w \models \varphi$  and  $\mathcal{M}|_\varphi, w \models \varphi$ . For  $\varphi$  to have a successful update in  $\mathcal{M}$ , there must be a disjunct  $\delta'$  of  $\varphi$  true at a point  $w$  in  $\mathcal{M}$  and a disjunct  $\delta$  of  $\varphi$  (possibly distinct from  $\delta'$ ) true at  $w$  in  $\mathcal{M}|_\varphi$ . What is required for  $\delta$  to be true at  $w$  in  $\mathcal{M}|_\varphi$  is that its propositional part  $\delta^\alpha$  is true at  $w$  in  $\mathcal{M}$ , that its universal part  $\delta^\square$  is already true at  $w$  in  $\mathcal{M}$  or *becomes* true at  $w$  in  $\mathcal{M}|_\varphi$ , and that its existential part  $\delta^\diamond$  remains true at  $w$  in  $\mathcal{M}|_\varphi$ . This is exactly what  $\chi$  captures.

**Proposition A.2** *Every strong Moore sentence is a Moore sentence.*

**Proof.** By Definitions A.1 and 3.12, it suffices to show that for  $\varphi$  in normal form, if a disjunct  $\delta$  of  $\varphi$  is not compensated in  $\varphi$ , then  $\delta$  is a Moore conjunction. To prove the contrapositive, suppose  $\delta$  is not a Moore conjunction. Then by Definition 3.7,  $\delta \wedge \diamond \delta^\alpha$  is clear and  $\delta \wedge \diamond (\delta^\alpha \wedge \gamma_k)$  is clear for every  $\diamond \gamma_k$  conjunct in  $\delta$ . It follows that  $\delta \wedge \bigwedge_{k \leq m} \diamond (\delta^\alpha \wedge \gamma_k)$  is clear and hence satisfiable. Given  $\mathcal{M}, w \models \delta \wedge \bigwedge_{k \leq m} \diamond (\delta^\alpha \wedge \gamma_k)$ , let  $S$  be the set of disjuncts in  $\varphi$  such that  $\mathcal{M}, w \models \bigwedge_{\sigma \in S} \sigma^{\square \diamond} \wedge \bigwedge_{\sigma \notin S} \neg \sigma^{\square \diamond}$ . For every  $\sigma \notin S$ , pick a conjunct  $\hat{\sigma}^{\square \diamond}$  of  $\sigma^{\square \diamond}$  such that  $\mathcal{M}, w \models \bigwedge_{\sigma \in S} \sigma^{\square \diamond} \wedge \bigwedge_{\sigma \notin S} \sim \hat{\sigma}^{\square \diamond}$ . Since  $\delta \in S$ , let the sequence  $\sigma_{\gamma_1}, \dots, \sigma_{\gamma_m}$  of disjuncts from  $S$  be such that  $\sigma_{\gamma_i} \equiv \delta$  for  $1 \leq i \leq m$ . We claim that for the  $\chi$  as in Definition A.1 based on these choices,  $\mathcal{M}, w \models \chi$ . From the fact that  $\mathcal{M}, w \models \bigwedge_{k \leq m} \diamond (\delta^\alpha \wedge \gamma_k)$  and our choice of  $\sigma_{\gamma_1}, \dots, \sigma_{\gamma_m}$ , we have  $\mathcal{M}, w \models \chi_\diamond$ . Since  $\mathcal{M}, w \models \delta^{\square}$ , it is immediate that  $\mathcal{M}, w \models \bigwedge_{\sigma \in S, j \leq n} \square (\sim \sigma \vee \beta_j)$ . Together with  $\mathcal{M}, w \models \bigwedge_{\sigma \in S} \sigma^{\square \diamond} \wedge \bigwedge_{\sigma \notin S} \sim \hat{\sigma}^{\square \diamond}$ , this gives  $\mathcal{M}, w \models \chi_\square$ . Finally, setting  $\delta' \equiv \delta$ , we have  $\mathcal{M}, w \models \chi$ . Hence  $\chi$  is clear, so  $\delta$  is compensated in  $\varphi$ .  $\square$

We can now generalize Theorem 3.13(i).

**Theorem A.3**  *$\varphi$  is self-refuting if and only if  $\varphi$  is a strong Moore sentence.*

**Proof.**  $\varphi$  is self-refuting iff any equivalent normal form of  $\varphi$  is self-refuting, so let us assume  $\varphi$  is already in normal form.

( $\Leftarrow$ ) We prove the contrapositive. Suppose  $\varphi$  is not self-refuting, so there is a pointed model such that  $\mathcal{M}, w \models \varphi$  and  $\mathcal{M}|_\varphi, w \models \varphi$ . Given  $\mathcal{M}|_\varphi, w \models \varphi$ , there is a disjunct  $\delta$  of  $\varphi$  such that  $\mathcal{M}|_\varphi, w \models \delta$ . We claim that  $\delta$  is compensated in  $\varphi$ , so  $\varphi$  is not a strong Moore sentence. It suffices to show the satisfiability of an appropriate  $\chi$  as in Definition A.1.

Since  $\mathcal{M}, w \models \varphi$ , there is a disjunct  $\delta'$  of  $\varphi$  such that  $\mathcal{M}, w \models \delta'$ . This gives the first conjunct of  $\chi$ . Since  $\mathcal{M}|_\varphi, w \models \delta$ , we have  $\mathcal{M}, w \models \delta^\alpha$  because  $\delta^\alpha$  is propositional and hence preserved under extensions. This gives the second conjunct of  $\chi$ .

Next we claim  $\mathcal{M}, w \models \chi_\square$ . Let  $S$  be the set of the disjuncts of  $\varphi$  such that  $\mathcal{M}, w \models \bigwedge_{\sigma \in S} \sigma^{\square \diamond} \wedge \bigwedge_{\sigma \notin S} \neg \sigma^{\square \diamond}$ , and let  $\hat{\sigma}^{\square \diamond}$  be a false conjunct of each  $\sigma^{\square \diamond}$  for  $\sigma \notin S$ . For *reductio*, suppose  $\mathcal{M}, w \not\models \bigwedge_{\sigma \in S, j \leq n} \square (\sim \sigma^\alpha \vee \beta_j)$ , where  $\square \beta_1 \wedge \dots \wedge \square \beta_n \equiv \delta^\square$ .

Then there is some  $v$  with  $wRv$  and  $\mathcal{M}, v \models \sigma^\alpha \wedge \neg \beta_j$  for some  $\sigma \in S$  and  $j \leq n$ . Since  $\sigma \in S$ , we have  $\mathcal{M}, w \models \sigma^{\square \diamond}$ . It follows by Lemma 3.8(i) that  $\mathcal{M}, v \models \sigma^{\square \diamond}$ , in which case  $\mathcal{M}, v \models \sigma$  given  $\mathcal{M}, v \models \sigma^\alpha$ . Hence  $v$  is retained in  $\mathcal{M}|_\varphi$ . But then given  $\mathcal{M}|_\varphi, v \not\models \beta_j$ , we have  $\mathcal{M}|_\varphi, w \not\models \square \beta_j$ , which contradicts the assumption that  $\mathcal{M}|_\varphi, w \models \delta$ . We conclude that  $\mathcal{M}, w \models \chi_\square$ , which gives the third conjunct of  $\chi$ .

Finally, we claim  $\mathcal{M}, w \models \chi_\diamond$ . Given  $\mathcal{M}|_\varphi, w \models \delta^\diamond$ , take an arbitrary  $\diamond \gamma_k$  in  $\delta$ , and let  $v$  be such that  $wRv$  and  $\mathcal{M}|_\varphi, v \models \gamma_k$ . It follows that  $\mathcal{M}, v \models \sigma$  for some disjunct  $\sigma$  of  $\varphi$ , which we label as  $\sigma_{\gamma_k}$ , for otherwise  $v$  would not be retained in  $\mathcal{M}|_\varphi$ . Since  $\gamma_k$  is propositional,  $\mathcal{M}, v \models \gamma_k$  given  $\mathcal{M}|_\varphi, v \models \gamma_k$ . Therefore  $\mathcal{M}, w \models \diamond (\sigma_{\gamma_k}^\alpha \wedge \gamma_k)$ . Then from the fact that  $\mathcal{M}, v \models \sigma_{\gamma_k}^{\square \diamond}$ , we have  $\mathcal{M}, w \models \sigma_{\gamma_k}^{\square \diamond}$  by Lemma 3.8(i), so  $\sigma_{\gamma_k} \in S$ . Since  $\diamond \gamma_k$  was arbitrary,  $\mathcal{M}, w \models \chi_\diamond$ , which gives the final conjunct of  $\chi$ .

( $\Rightarrow$ ) Again we prove the contrapositive. Suppose  $\delta$  is not a strong Moore sen-

tence, so there is some  $\delta$  that is compensated in  $\varphi$ , for which an appropriate  $\chi$  as in Definition A.1 is clear. Where  $\mathcal{M}, w \models \chi$ , we claim that  $\mathcal{M}_{|\varphi}, w \models \delta$ . We will show  $\mathcal{M}_{|\varphi}, w \models \delta^\alpha$ ,  $\mathcal{M}_{|\varphi}, w \models \delta^\square$ , and  $\mathcal{M}_{|\varphi}, w \models \delta^\diamond$  separately.

Given  $\mathcal{M}, w \models \delta$ , we have  $\mathcal{M}_{|\varphi}, w \models \delta^\alpha$  since  $\delta^\alpha$  is propositional.

Next, for any  $v$  retained in  $\mathcal{M}_{|\varphi}$ , we have  $\mathcal{M}, v \models \sigma^\alpha$  for some disjunct  $\sigma$  of  $\varphi$ . It must be that  $\sigma \in S$ , for otherwise  $\mathcal{M} \models \neg\sigma$  given  $\mathcal{M}, w \models \chi_\square$  and Lemma 3.8(i). Then from the fact that  $\mathcal{M}, w \models \bigwedge_{\sigma \in S, j \leq n} \square(\sim \sigma^\alpha \vee \beta_j)$ , we have  $\mathcal{M}_{|\varphi}, v \models \beta_j$  for all  $j \leq n$ . Since  $v$  was arbitrary,  $\mathcal{M}_{|\varphi}, w \models \delta^\square$ .

Finally, given  $\mathcal{M}, w \models \chi_\diamond$ , for any  $\diamond\gamma_k$  in  $\delta$  we have  $\mathcal{M}, w \models \diamond(\sigma_{\gamma_k} \wedge \gamma_k)$  with  $\sigma_{\gamma_k} \in S$ . Let  $v$  be such that  $wRv$  and  $\mathcal{M}, v \models \sigma_{\gamma_k}^\alpha \wedge \gamma_k$ . Since  $\sigma_{\gamma_k} \in S$ ,  $\mathcal{M}, w \models \sigma_{\gamma_k}^{\square\diamond}$ . It follows by Lemma 3.8(i) that  $\mathcal{M}, v \models \sigma_{\gamma_k}^{\square\diamond}$ , in which case  $\mathcal{M}, v \models \sigma_{\gamma_k}$  given  $\mathcal{M}, v \models \sigma_{\gamma_k}^\alpha$ . Hence  $v$  is retained in  $\mathcal{M}_{|\varphi}$ . Since  $\gamma_k$  is propositional, we have  $\mathcal{M}_{|\varphi}, v \models \gamma_k$  given  $\mathcal{M}, v \models \gamma_k$ , whence  $\mathcal{M}_{|\varphi}, w \models \diamond\gamma_k$ . Since  $\gamma_k$  was arbitrary,  $\mathcal{M}_{|\varphi}, w \models \delta^\diamond$ .

We conclude that  $\mathcal{M}_{|\varphi}, w \models \delta$ , in which case  $\mathcal{M}_{|\varphi}, w \models \varphi$ . Given  $\mathcal{M}, w \models \chi$ , we have  $\mathcal{M}, w \models \delta$  and hence  $\mathcal{M}, w \models \varphi$ , so  $\varphi$  is not self-refuting.  $\square$

Finally, we will prove an analogous generalization of Theorem 3.13(ii).

**Definition A.4**  $\varphi$  is a *strong Moorean sentence* iff for any normal form  $\varphi^*$  of  $\varphi$ , there is a disjunct  $\delta$  and non-empty sets  $S$  and  $T$  of disjuncts of  $\varphi^*$ , with for every  $\theta \in T$ , a  $\diamond\gamma_\theta$  in  $\theta$ , such that  $\chi \equiv \delta \wedge \chi_1 \wedge \chi_2 \wedge \chi_3$  is clear, where

$$\begin{aligned} t(\theta) &\equiv \theta^{\alpha\diamond} \wedge \bigwedge_{\sigma \in S, \square\beta \text{ in } \theta} \square(\sim \sigma^\alpha \vee \beta); \\ \chi_1 &\equiv \bigwedge_{\theta \in T} t(\theta) \wedge \bigwedge_{\theta \notin T} \neg t(\theta); \quad \chi_2 \equiv \bigwedge_{\sigma \in S} \sigma^{\square\diamond} \wedge \bigwedge_{\sigma \notin S} \neg \sigma^{\square\diamond}; \\ \chi_3 &\equiv \bigwedge_{\sigma \in S, \theta \in T} \square(\sim \sigma^\alpha \vee \sim \gamma_\theta). \end{aligned}$$

The formula  $\chi$  is not yet in normal form, due to the  $\neg t(\theta)$  conjuncts in  $\chi_1$  and the  $\neg \sigma^{\square\diamond}$  conjuncts in  $\chi_2$ , so strictly the definition of clarity does not apply. However, it is straightforward to put  $\chi$  into normal form using the same method involving  $\hat{\sigma}^{\square\diamond}$  as in Definition A.1, together with some distribution of  $\wedge$  and  $\vee$ . Since in this case the necessary modifications add four existential quantifiers to the definition, for simplicity we do not write them out. When we say that  $\chi$  is clear, strictly we mean that the modified formula is clear.

As we will see in the proof of Theorem A.6, the clarity of  $\chi$  is equivalent to there being a model in which  $\varphi$  has an unsuccessful update. For  $\varphi$  to have an unsuccessful update in  $\mathcal{M}$ , there must be some disjunct  $\delta$  in  $\varphi$  such that  $\mathcal{M}, w \models \delta$ , but no disjunct  $\delta'$  such that  $\mathcal{M}_{|\varphi}, w \models \delta'$ . To ensure that there are no such  $\delta'$ , we need only keep track of those disjuncts  $\theta$  of  $\varphi$  whose propositional part  $\theta^\alpha$  and existential part  $\theta^\diamond$  are true in  $\mathcal{M}$  and whose universal part  $\theta^\square$  was already true in  $\mathcal{M}$  or *becomes* true in  $\mathcal{M}_{|\varphi}$ , since all other disjuncts will be false in  $\mathcal{M}_{|\varphi}$ . This is the purpose of  $\chi_1$ . For each such  $\theta$ , we must ensure that there is a diamond formula  $\diamond\gamma_\theta$  in  $\theta$  that becomes false in  $\mathcal{M}_{|\varphi}$ , because none of its witnesses satisfy any of the disjuncts that are satisfied somewhere in  $\mathcal{M}$ . This is the purpose of  $\chi_3$  and  $\chi_2$ .

**Proposition A.5** *Every strong Moorean sentence is a Moorean sentence.*

**Proof.** Assume  $\varphi$  is a strong Moorean sentence, so an appropriate  $\chi$  as in Definition A.4 is clear. Note that the distinguished disjunct  $\delta$  of  $\varphi^*$  must be a member of both  $S$  and  $T$ . Then given that  $\delta \wedge \chi_3$  is clear, we have that  $\delta \wedge \Box(\sim \delta^\alpha \vee \sim \gamma_\delta)$  is clear, where  $\Diamond\gamma_\delta$  is a conjunct of  $\delta$ . Hence  $\delta$  is a Moorean conjunction by Definition 3.7, in which case  $\varphi$  is a Moorean sentence by Definition 3.12.  $\square$

**Theorem A.6**  *$\varphi$  is unsuccessful if and only if  $\varphi$  is a strong Moorean sentence.*

**Proof.** As before, let us assume that  $\varphi$  is already in normal form.

( $\Leftarrow$ ) Suppose  $\varphi$  is a strong Moorean sentence, so an appropriate  $\chi$  as in Definition A.4 is clear. Consider a pointed model such that  $\mathcal{M}, w \vDash \chi$ . For *reductio*, assume  $\mathcal{M}_{|\varphi}, w \vDash \varphi$ . Then there exists a disjunct  $\theta$  in  $\varphi$  such that  $\mathcal{M}_{|\varphi}, w \vDash \theta$ .

We claim that  $\theta \in T$ . For if  $\theta \notin T$ , then given  $\mathcal{M}, w \vDash \chi_1$  there are two cases. Case 1:  $\mathcal{M}, w \not\vDash \theta^{\alpha\Diamond}$ . Then  $\mathcal{M}_{|\varphi}, w \not\vDash \theta^{\alpha\Diamond}$  since  $\theta^{\alpha\Diamond}$  is existential. Case 2:  $\mathcal{M}, w \vDash \Diamond(\sigma^\alpha \wedge \neg\beta)$  for some  $\sigma \in S$  and  $\Box\beta$  in  $\theta$ . Let  $v$  be such that  $wRv$  and  $\mathcal{M}, v \vDash \sigma^\alpha \wedge \neg\beta$ . Since  $\sigma \in S$ , we have  $\mathcal{M}, w \vDash \sigma^{\Box\Diamond}$ , in which case  $\mathcal{M}, v \vDash \sigma$  given Lemma 3.8(i) and  $\mathcal{M}, v \vDash \sigma^\alpha$ . Hence  $v$  is retained in  $\mathcal{M}_{|\varphi}$ , and since  $\beta$  is propositional,  $\mathcal{M}_{|\varphi}, v \not\vDash \beta$ , so  $\mathcal{M}_{|\varphi}, w \not\vDash \Box\beta$  and  $\mathcal{M}_{|\varphi}, w \not\vDash \theta^\square$ . In both cases,  $\mathcal{M}_{|\varphi}, w \not\vDash \theta$ , a contradiction. Therefore  $\theta \in T$ .

Given  $\mathcal{M}_{|\varphi}, w \vDash \theta$ , for every  $\Diamond\gamma$  in  $\theta$  there is a  $v$  with  $wRv$  and  $\mathcal{M}_{|\varphi}, v \vDash \gamma$ . Since  $v$  was retained in  $\mathcal{M}_{|\varphi}$ , we have  $\mathcal{M}, v \vDash \sigma$  for some  $\sigma \in S$ . Then given  $\mathcal{M}, w \vDash \chi_3$  and  $\theta \in T$ , we have  $\mathcal{M}, v \vDash \neg\gamma$ . Since  $\gamma$  is propositional,  $\mathcal{M}_{|\varphi}, v \vDash \neg\gamma$ , a contradiction. We conclude that  $\mathcal{M}_{|\varphi}, w \not\vDash \varphi$ , so  $\varphi$  is unsuccessful.

( $\Rightarrow$ ) Suppose  $\varphi$  is unsuccessful, so there is a pointed model with  $\mathcal{M}, w \vDash \varphi$  but  $\mathcal{M}_{|\varphi}, w \not\vDash \varphi$ . To show that  $\varphi$  is a strong Moorean sentence, it suffices to show that an appropriate  $\chi$  as in Definition A.4 is satisfiable. Given  $\mathcal{M}, w \vDash \varphi$ , we can read off from  $w$  the disjunct  $\delta$  and sets  $S$  and  $T$  such that  $\mathcal{M}, w \vDash \delta \wedge \chi_1 \wedge \chi_2$ . It only remains to show  $\mathcal{M}, w \vDash \chi_3$ . For *reductio*, suppose there is a  $\theta \in T$  such that for all  $\Diamond\gamma$  in  $\theta$ ,  $\mathcal{M}, w \not\vDash \bigwedge_{\sigma \in S} \Box(\sim \sigma^\alpha \vee \sim \gamma)$ , i.e.,  $\mathcal{M}, w \vDash \Diamond(\sigma^\alpha \wedge \gamma)$  for some  $\sigma \in S$ . Then we claim that  $\mathcal{M}_{|\varphi}, w \vDash \theta$ . As before, we take the three parts of  $\theta$  separately.

For  $\theta^\Diamond$ , consider any  $\Diamond\gamma$  in  $\theta$  and let  $v$  be such that  $wRv$  and  $\mathcal{M}, v \vDash \sigma^\alpha \wedge \gamma$ . Given  $\mathcal{M}, w \vDash \chi_2$  and the fact that  $\sigma \in S$ , we have  $\mathcal{M}, w \vDash \sigma^{\Box\Diamond}$ , so  $\mathcal{M}, v \vDash \sigma$  by Lemma 3.8(i). Hence  $v$  is retained in  $\mathcal{M}_{|\varphi}$ . Since  $\mathcal{M}, v \vDash \gamma$  and  $\gamma$  is propositional,  $\mathcal{M}_{|\varphi}, v \vDash \gamma$  and therefore  $\mathcal{M}_{|\varphi}, w \vDash \Diamond\gamma$ . Since  $\Diamond\gamma$  was arbitrary,  $\mathcal{M}_{|\varphi}, w \vDash \theta^\Diamond$ .

For  $\theta^\square$ , take any  $v$  retained in  $\mathcal{M}_{|\varphi}$  such that  $wRv$ , and we have  $\mathcal{M}, v \vDash \sigma^\alpha$  for some  $\sigma \in S$ . It follows that for any  $\Box\beta$  in  $\theta$ , we have  $\mathcal{M}, v \vDash \beta$  given  $\mathcal{M}, w \vDash \chi_1$ . Since  $\beta$  is propositional,  $\mathcal{M}_{|\varphi}, v \vDash \beta$ . Since  $v$  and  $\beta$  were arbitrary,  $\mathcal{M}_{|\varphi}, w \vDash \delta^\square$ .

Finally, for  $\theta^\alpha$ , since by assumption  $\theta \in T$  and  $\mathcal{M}, w \vDash \chi_1$ , we have  $\mathcal{M}, w \vDash \theta^\alpha$  and hence  $\mathcal{M}_{|\varphi}, w \vDash \theta^\alpha$ .

We have shown  $\mathcal{M}_{|\varphi}, w \vDash \theta$  and hence  $\mathcal{M}_{|\varphi}, w \vDash \varphi$ , which contradicts our initial assumption. It follows that for every  $\theta \in T$  there is a  $\Diamond\gamma_\theta$  in  $\theta$  such that  $\mathcal{M}, w \vDash \chi_3$ .  $\square$