
A Foundation for Metareasoning Part II: The Model Theory

GIOVANNI CRISCUOLO, *Dipartimento di Scienze Fisiche, University of Naples, 80125 Napoli, Italy.*
E-mail: vanni@na.infn.it

FAUSTO GIUNCHIGLIA, *Dipartimento di Informatica e Studi Aziendali, University of Trento, 38100 Trento, Italy, and ITC-IRST, Centro per la Ricerca Scientifica e Tecnologica, 38050 Trento, Italy.*
E-mail: fausto@itc.it

LUCIANO SERAFINI, *ITC-IRST, Centro per la Ricerca Scientifica e Tecnologica, 38050 Trento, Italy.*
E-mail: serafini@itc.it

Abstract

OM pairs are our proposed framework for the formalization of metareasoning. OM pairs allow us to generate deductively the object theory and/or the metatheory. This is done by imposing, via appropriate reflection rules, the relation we want to hold between the object theory and the metatheory. In a previous paper we have studied the proof theoretic properties of OM pairs. In this paper we study their model theoretic properties, in particular we study the relation between the models of the metatheory and the object theory; and how to use these results to refine the previous analysis.

Keywords: Metatheory, reflection, contextual reasoning.

1 Introduction

The '*Meta*' property is not a property which can be ascribed to a single theory; it is rather a relation between two theories. A theory is said to be *meta* of another, not necessarily distinct, theory if it is *about* this other theory. The latter theory is often said to be the *object* of the former theory. The metatheory speaks about the object theory by using a special set of predicates, called *metapredicates*, each of which is supposed to represent a property about the object theory. By property here we mean a set of formulas of the language of object theory. Examples of properties are: 'being true in a model of the object theory', 'being a theorem in the object theory', 'being a theorem in the object theory extended with the set of axioms X' etc.

A *reflection principle* [3] expresses the correctness of the metatheory when a metapredicate is interpreted in a certain property about the object theory. A reflection principle, therefore, is a statement of the form: '*if a metaformula, expressing a certain proposition about the object theory, is a theorem of the metatheory, then such a proposition is actually true*'. The definition of metatheories that satisfy a reflection principle has been one of the main research issues in the area of Formal Logic, Computer Science, and Artificial Intelligence. Such metathe-

ories are usually built by adding to a theory special formulas, or special inference rules, or combinations of the two for the metapredicates.

In a Part I [2] we introduced a new framework called OM pairs for a uniform specification of metatheories that satisfy different reflection principles. OM pairs allow us to specify metatheories via a set of inference rules, called *reflection rules*. Reflection rules are special inference rules between the object theory and the metatheory (which are supposed to be distinct theories). They allow one to infer formulas in the metatheory starting from formulas in the object theory and vice versa.

Once a metatheory is specified in the previous terms, it is necessary to characterize it, namely, to verify whether it satisfies a given reflection principle. This amounts to verifying, whether the metapredicate correctly represents the property expressed in the reflection principle. This characterization concerns two aspects: the set of theorems of the metatheory, and its class of models. The former problem has been studied in Part I [2]. In this paper we concentrate on the latter. We propose a model theoretic characterization of the metatheory in terms of

the class of properties about the object theory (represented as set of formulas) which are the interpretations of the metapredicate in a model of the metatheory. (1.1)

Notice that we consider a class of properties rather than a single one. This is because, in general, there is more than one interpretation of the metapredicate that satisfies the metatheory. (1.1) can be rephrased more formally as follows:

the class of sets of formulas P in the language of the object theory, such that there is a model m_P of the metatheory that interprets the metapredicate in P . (1.2)

The core of the paper contains a characterization of metatheories containing a single metapredicate \bullet , and specified via OM pairs. We show how metatheories which respect reflection principles about different properties of the object theory, such as ‘the set of formulas satisfiable in a specific model of the object theory’, ‘the set of formulas that are provable in an extension of the object theory’ etc. can be defined via different combinations of reflection rules.

The model-theoretic characterization allows us to refine the characterization of the different reflection rules given in Part I [2]. We indeed refine the partial order on the strength of reflection rules given in that paper, by proving that it is indeed a strict partial order. The consequence of this is that each combination of reflection rules represents different reflection principles.

The paper is structured as follows: in Section 2 we recall the basic definition of OM pair given in Part I [2]. In Section 3 we characterize the models of the metatheory. The main idea is to define them as sets of object level formulas. This allows us to study the relations between the models of the metatheory and the object theory and, therefore, to make precise statements like ‘the metatheory speaks about the object theory’ and ‘the object theory is the model of the metatheory’. Section 4 refines the analysis described in Section 3 of [2] and gives a strict ordering characterization of the reflection principles formalizables by OM pairs, based on the strength of the derivability relations and the sets of theorems generated. In Section 5, we generalize the notion of duality, given in [2], to models, and we show that it preserves satisfiability.

2 OM pairs

In this section we briefly recall the main definition of OM pairs given in Part I.

We start from two arbitrary theories O and M , presented as axiomatic formal systems, i.e. $O = \langle L_O, \Omega_O, \Delta_O \rangle$ and $M = \langle L_M, \Omega_M, \Delta_M \rangle$. O is called the *object theory*, M is called the *metatheory*. L_O (L_M) is the *language* of O (M), Ω_O (Ω_M) is the *set of axioms* of O (M), and Δ_O (Δ_M) is the *deductive machinery* (set of inference rules) of O (M). We call *Meta-formulas*, the formulas of L_M . We suppose that Δ_O and Δ_M are set of natural deduction style inference rules, as defined in [6].

\vdash_O and \vdash_M , are the derivability relations defined by O and M , respectively; $\text{TH}(O)$ and $\text{TH}(M)$ are the set of theorems of O and M respectively. We then consider new kinds of inter-theory inference rules, called *bridge rules* [4], whose premisses and conclusions belong to different languages. Thus we may have a bridge rule with premisses in L_O and conclusion in L_M ; and, vice versa, a bridge rule with premisses in L_M and conclusion in L_O . Bridge rules extend deductively $\text{TH}(O)$ and $\text{TH}(M)$, that is, new object and metatheorems may be proved by applying bridge rules. The set of theorems of the metatheory may be extended whenever the set of theorems of the object theory is extended, and vice versa. We focus on the following bridge rules, and we call them *reflection rules*.

$$\frac{A}{\bullet("A")} \text{Rup} \quad \frac{\neg A}{\neg \bullet("A")} \text{Rup}^n \quad \frac{\bullet("A")}{A} \text{Rdw} \quad \frac{\neg \bullet("A")}{\neg A} \text{Rdw}^n$$

For each of the reflection rule above, we consider two versions, the unrestricted version, which can be applied with no restriction and the restricted version (denoted by adding the index), which can be applied under the following condition:

RESTRICTIONS: Rules labelled with index r are applicable if the premiss does not depend on any assumptions in the same theory.

For the rest of the paper we make the following simplifying hypotheses: L_O and L_M are propositional and L_M is the *propositional metalanguage* of L_O .¹

DEFINITION 2.1 (Propositional metalanguage)

Given a logical language L , its propositional metalanguage is the propositional language $\bullet("L")$, whose set of atomic wffs is the set $\{\bullet("A") : A \in L\}$.

DEFINITION 2.2 (OM pair)

An *Object-Meta Pair* (OM pair) is a triple $\langle O, M, (\text{RR}) \rangle$ where $O = \langle L_O, \Omega_O, \Delta_O \rangle$, $M = \langle L_M, \Omega_M, \Delta_M \rangle$, L_M contains the propositional metalanguage of L_O , and (RR) is a set of reflection rules.

Given an object theory O , a metatheory M , and a set of reflection rules (RR) , we say that $\text{OM} = \langle O, M, (\text{RR}) \rangle$ is the OM pair *composed of O and M connected by (RR)* . Notationally, when it contains more than one bridge rule we represent (RR) by listing its elements separated by a $+$. One example of (RR) is $\text{Rup} + \text{Rdw}^n$.

A wff $A \in L_O \cup L_M$ is derivable from a set of assumptions $\Gamma \subseteq L_O \cup L_M$, in OM, in symbols $\Gamma \vdash_{\text{OM}} A$, if there is a deduction in OM of A from Γ that applies Δ_O , Δ_M , the axioms of O and M , and the reflection rules of OM. A is *provable in (a theorem of) OM*, abbreviated as $\vdash_{\text{OM}} A$, if it is derivable from the empty set.

$\text{TH}_{\text{OM}}(O) \subseteq L_O$ is the set of formulas of L_O provable in OM. $\text{TH}_{\text{OM}}(M) \subseteq L_M$ is the set of formulas of L_M provable in OM. Trivially $\text{TH}(O) \subseteq \text{TH}_{\text{OM}}(O)$ and $\text{TH}(M) \subseteq$

¹In Part I [2] we have shown how the results in the propositional case can be generalized to the first-order case.

$\text{TH}_{\text{OM}}(M)$. Given an OM pair $\text{OM} = \langle O, M, (\text{RR}) \rangle$ we use the term *object theory of OM* (generated by (RR)) to denote any theory $O' = \langle L_O, \Omega'_O, \Delta_O \rangle$ whose set of theorems is $\text{TH}_{\text{OM}}(O')$ and the term *metatheory of OM* (generated by (RR)) to denote any theory $M' = \langle L_M, \Omega'_M, \Delta_M \rangle$ whose set of theorems is $\text{TH}_{\text{OM}}(M')$. Trivially $\text{TH}(O') = \text{TH}_{\text{OM}}(O)$ and $\text{TH}(M') = \text{TH}_{\text{OM}}(M)$.

We consider the OM pairs with the following combinations of reflection rules:

- | | |
|---|---|
| 1. \emptyset | 12. $\text{Rdw} + \text{Rup}_r^n + \text{Rdw}_r^n$ |
| 2. Rdw_r | 13. $\text{Rup}_r + \text{Rdw}$ |
| 3. Rup_r | 14. $\text{Rup}_r + \text{Rdw} + \text{Rup}_r^n$ |
| 4. Rdw | 15. $\text{Rup}_r + \text{Rdw} + \text{Rdw}_r^n$ |
| 5. Rup | 16. $\text{Rup}_r + \text{Rdw} + \text{Rup}^n$ |
| 6. $\text{Rup}_r + \text{Rdw}_r$ | 17. $\text{Rup}_r + \text{Rdw} + \text{Rdw}^n$ |
| 7. $\text{Rup}_r + \text{Rdw}_r^n$ | 18. $\text{Rup}_r + \text{Rdw} + \text{Rup}_r^n + \text{Rdw}_r^n$ |
| 8. $\text{Rup}_r + \text{Rdw}_r + \text{Rdw}_r^n$ | 19. $\text{Rup}_r + \text{Rdw} + \text{Rup}_r^n + \text{Rdw}^n$ |
| 9. $\text{Rup}_r + \text{Rdw}_r + \text{Rup}_r^n$ | 20. $\text{Rup} + \text{Rdw}$ |
| 10. $\text{Rup}_r + \text{Rdw}_r + \text{Rup}_r^n + \text{Rdw}_r^n$ | 21. $\text{Rup} + \text{Rdw} + \text{Rup}_r^n$ |
| 11. $\text{Rdw} + \text{Rup}_r^n$ | 22. $\text{Rup} + \text{Rdw}^n$ |

3 A model-theoretic characterization

The goal of this section is to understand the relation between the models of the metatheory and the object theory. For our purpose it is enough to consider the case in which L_M is the propositional metalanguage of L_O , i.e. all its atomic wffs are of the form $\bullet(\text{"}A\text{"})$ for some wff A of L_O . Being a propositional theory, the models of M are truth assignments for the atomic wffs of L_M . Any interpretation m of L_M identifies a subset $S_m \subseteq L_O$ according to the following definition:

$$S_m = \{A \in L_O \mid m(\bullet(\text{"}A\text{"})) = \text{True}\}.$$

We can therefore consider the set of interpretations of L_M to be the powerset of the wffs in L_O .

For any set (RR) of reflection rules, we characterize:

- the models of the metatheory generated by (RR) in terms of O ;
- the set of theorems of the object theory generated by (RR) in terms of the models of M .

The first result allows us to make precise in which sense *the metatheory speaks about the object theory*. The second result allows us to make precise the sense in which *the object theory is a model of the metatheory*.

We use the following notation. Let m be an interpretation of L_M , we say that m is a model of Γ , in symbols $m \models \Gamma$, to mean that m satisfies all wffs in Γ (according to the classical definition of satisfiability).

Furthermore we say that a set of wffs $\Gamma \subseteq L$ is a *classical propositional theory* if it contains all the classical tautologies and it is closed under modus ponens. Γ is consistent if $\perp \notin \Gamma$. Γ is *maximal* if for any wff $A \in L$, $A \in \Gamma$ or $\neg A \in \Gamma$. We also use the following notation. For any theory $T = \langle L, \Omega, \Delta \rangle$ and any $\Gamma \subseteq L$, $T + \Gamma$ denotes the extended theory $\langle L, \Omega \cup \Gamma, \Delta \rangle$. For any set of wffs $\Gamma \subseteq L_O$, $\text{Rup}(\Gamma)$ denotes the set of wffs $\{\bullet(\text{"}A\text{"}) : A \in \Gamma\} \subseteq L_M$. We sometimes use $\bullet(\text{"}\Gamma\text{"})$ instead of $\text{Rup}(\Gamma)$. Analogously for any set of wffs $\Gamma \subseteq L_M$, $\text{Rdw}(\Gamma)$ denotes the set of wffs $\{A : \bullet(\text{"}A\text{"}) \in \Gamma\} \subseteq L_O$.

3.1 The models of the metatheory

Let us start by considering the metatheories generated by combinations of reflection rules which are characterizable axiomatically, i.e. the metatheories for which $\text{TH}_{\text{OM}}(M) = \text{TH}(M + \Phi)$, with Φ being some schematic metaformula. ([2, Section 4] gives an in depth analysis of the axiomatic characterizability of metatheories.) In this case the (obvious) idea is to look for the subclass of interpretations of L_M which satisfies Ω_M and Φ .

THEOREM 3.1

Let OM be an OM pair composed of O and M connected by the set of reflection rules (RR);

1. if $(\text{RR}) =_M \text{Rdw}_r$, then m is a model of $\text{TH}_{\text{OM}}(M)$, if and only if it is a model of Ω_M ;
2. if $(\text{RR}) =_M \text{Rup}_r$, then m is a model of $\text{TH}_{\text{OM}}(M)$, if and only if it is a model of Ω_M , and $\text{TH}(O) \subseteq m$.

Item (i) of Theorem 3.1 states that any model of M is a model of the metatheory of OM and vice versa. m is a model of the metatheory of OM independently of the object theory. Item (ii) states that the models of the metatheory of OM are all the models of Ω_M which contain the theorems of the object theory.

PROOF. [Theorem 3.1] These results are a direct consequence of [2, Theorem 4.1]. This theorem states the (obvious) result that reflection down leaves the metatheory unchanged, and we have $\text{TH}_{\text{OM}}(M) = \text{TH}(M)$, while reflection up adds to the metatheory all the theorems of the form $\bullet("A")$, with $A \in \text{TH}(O) = \text{TH}_{\text{OM}}(O)$. ■

THEOREM 3.2

Let OM be an OM pair composed of O and M connected by the set of reflection rules (RR). An interpretation m of L_M is a model of $\text{TH}_{\text{OM}}(M)$ if and only if, $m \models \Omega_M$ and:

1. If (RR) is $\text{Rup}_r + \text{Rdw}$ then, $\text{TH}(O) \subseteq m$ and m is a *classical propositional theory*.
2. If (RR) is $\text{Rup}_r + \text{Rdw} + \text{Rup}_r^n$, then $\text{TH}(O) \subseteq m$ and m is a *consistent classical propositional theory*.
3. If (RR) is $\text{Rup} + \text{Rdw}$, then $\text{TH}(O) \subseteq m$ and m is a *maximal classical propositional theory*.
4. If (RR) is $\text{Rup} + \text{Rdw} + \text{Rup}_r^n$, then $\text{TH}(O) \subseteq m$ and m is a *consistent and maximal classical propositional theory*.
5. If (RR) is either $\text{Rup}_r^n + \text{Rdw}$ or $\text{Rup}_r^n + \text{Rdw}$, then either $O + m$ is *consistent*, or $m = \emptyset$.

To understand Theorem 3.2, let us consider the case with $\Omega_M = \emptyset$. In this situation the models of $\text{TH}_{\text{OM}}(M)$ are completely determined by O and (RR). (These observations can be easily generalized. In fact, from Theorem 3.2, $\Omega_M \neq \emptyset$ restricts us to the models of Ω_M .) Consider item (1). The models of $\text{TH}_{\text{OM}}(M)$ are the sets of theorems of the extensions of the object theory, i.e. all the theories which can be obtained by adding a set of axioms Γ to O . Therefore, if A is a logical consequence of a set of formulas Γ , then $\bullet("A")$ is also a logical consequence of $\bullet(" \Gamma ")$. Consider item (2). The models of $\text{TH}_{\text{OM}}(M)$ are the sets of theorems of the consistent extensions of the object theory. Notice that, as there are no consistent extensions of an inconsistent theory, if O is inconsistent then $\text{TH}_{\text{OM}}(M)$ is inconsistent as well. Consider item (3). The models of $\text{TH}_{\text{OM}}(M)$ are the maximal (possibly inconsistent) extensions of the object theory. Consider item (4). The models of $\text{TH}_{\text{OM}}(M)$ are the

maximal consistent extensions of the object theory. Since the set of maximal consistent sets containing the object theory is isomorphic to the set of models of the object theory, $\bullet("A")$ can be interpreted as ‘ A is true in m ’. Finally, consider item (5). The models of $\text{TH}_{\text{OM}}(M)$ are the sets of formulas which can be added to the object theory without making it inconsistent. If, for instance, A is derivable from Γ in O , the object theory extended by Γ and $\neg A$ is inconsistent. This corresponds to the fact that in any model m of $\text{TH}_{\text{OM}}(M)$ it is never the case that Γ and $\neg A$ both belong to m . A further example: if the object theory is inconsistent (i.e. all formulas are provable) then it cannot be consistently extended. In this case the only model of the metatheory of OM is the empty set.

PROOF. [Theorem 3.2] This proof is based on the result of Theorems 4.2 and 4.3 in [2]. In these theorems the metatheory of OM is characterized in terms of the following formulas:

$$\begin{aligned}
& \bullet("A \supset B") \supset (\bullet("A") \supset \bullet("B")) && \text{(K)} \\
& \neg \bullet("A") \supset \bullet("\neg A") && \text{(Comp)} \\
& \neg \bullet("\perp") && \text{(nTbot)} \\
& \{\bullet("A_1") \wedge \dots \wedge \bullet("A_k") \supset \neg \bullet("A")\} \vdash_{O+A_1, \dots, A_n} \neg A && \text{(Cons)}
\end{aligned}$$

where A, B are understood as schematic variables (parameters) ranging over L_O .

Item (1). In this case $\text{TH}_{\text{OM}}(M) = \text{TH}(M + \text{Rup}(\text{TH}(O)) + (K))$, as proved by [2, Theorem 4.2]. Then, $m \models \text{TH}_{\text{OM}}(M)$ if and only if $m \models \Omega_M$, $m \models \text{Rup}(\text{TH}(O))$ and $m \models (K)$. Suppose that m is a model of $\text{TH}_{\text{OM}}(M)$, then $m \models \Omega_M$ and $\text{TH}(O) \subseteq m$. m is a classical propositional theory as it contains all the propositional tautologies (contained in $\text{TH}(O)$), and it is closed under *modus ponens* since $m \models (K)$. Vice versa, suppose that $m \models \Omega_M$ and m is a classical propositional theory which contains $\text{TH}(O)$, then $m \models \Omega_M$, by hypothesis, $m \models \text{Rup}(\text{TH}(O))$ as $\text{TH}(O) \subseteq m$, and $m \models (K)$ as m is closed under logical consequence.

Item (2). In this case $\text{TH}_{\text{OM}}(M) = \text{TH}(M + \text{Rup}(\text{TH}(O)) + (K) + (\text{nTbot}))$, as proved by [2, Theorem 4.2]. Then, $m \models \text{TH}_{\text{OM}}(M)$ if and only if, $m \models \Omega_M$, $m \models \text{Rup}(\text{TH}(O))$, $m \models (K)$ and $m \models (\text{nTbot})$. From item (1), if $m \models \text{TH}_{\text{OM}}(M)$, then $m \models \Omega_M$, and m is a classical propositional theory that contains $\text{TH}(O)$. Furthermore $m \models (\text{nTbot})$ implies that $\perp \notin m$ and therefore that m is consistent. Vice versa, if $m \models \Omega_M$ and m is a consistent classical propositional theory which contains $\text{TH}(O)$ then, by item (1), $m \models \text{Rup}(\text{TH}(O))$ and $m \models (K)$. The consistency of m implies that $\perp \notin m$ and therefore that $m \models (\text{nTbot})$.

Item (3). In this case $\text{TH}_{\text{OM}}(M) = \text{TH}(M + \text{Rup}(\text{TH}(O)) + (K) + (\text{Comp}))$, as proved by [2, Theorem 4.2]. Then, $m \models \text{TH}_{\text{OM}}(M)$ if and only if $m \models \Omega_M$, $m \models \text{Rup}(\text{TH}(O))$, $m \models (K)$ and $m \models (\text{Comp})$. From item (1) of this theorem, if $m \models \text{TH}_{\text{OM}}(M)$, then $m \models \Omega_M$, and m is a classical propositional theory that contains $\text{TH}(O)$. Furthermore $m \models (\text{Comp})$ implies that, for any object wff A , $A \in m$ or $\neg A \in m$, i.e. m is maximal. Vice versa, if $m \models \Omega_M$ and m is a maximal classical propositional theory which contains $\text{TH}(O)$, then, by item (1) of this theorem, $m \models \text{Rup}(\text{TH}(O))$ and $m \models (K)$. Furthermore m being maximal implies that $A \in m$ or $\neg A \in m$ and therefore that $m \models (\text{Comp})$.

Item (4). In this case $\text{TH}_{\text{OM}}(M) = \text{TH}(M + \text{Rup}(\text{TH}(O)) + (K) + (\text{Comp}) + (\text{nTbot}))$, as proved by [2, Theorem 4.2]. Item (4) is therefore a consequence of items (2) and (3) of this theorem.

Item (5). In this case $\text{TH}_{\text{OM}}(M) = \text{TH}(M + (\text{Cons}))$, as proved by [2, Theorem 4.3]. Let us start with the only if direction. Let m be a model of $\text{TH}_{\text{OM}}(M)$. Then, m is a model

of Ω_M , and by [2, Theorem 4.3], for any A_1, \dots, A_n, A such that $\vdash_{O+A_1, \dots, A_n} \neg A$

if $A_1, \dots, A_n \in m$, then $A \notin m$.

Suppose that $O + m$ is inconsistent. Then, for any formula A , $\vdash_{O+m} \neg A$, this implies that there is a finite number of wffs $A_1, \dots, A_n \in m$ such that $\vdash_{O+A_1, \dots, A_n} \neg A$. This implies that $A \notin m$. As this is true for any formula A , m must be the empty set.

As far as the only if direction, to prove that $m \models \text{TH}_{\text{OM}}(M)$ we exploit again [2, Theorem 4.3] and prove that

$$m \models M + \{\bullet("A_1") \wedge \dots \wedge \bullet("A_n") \supset \neg \bullet("A") : \vdash_{O+A_1, \dots, A_n} \neg A\}. \quad (3.1)$$

If $m = \emptyset$ then $m \models \neg \bullet("A")$ for any formula A , and since m satisfies Ω_M , then (3.1) is true. Suppose instead $m \neq \emptyset$ and $m \models \Omega_M$. Then, suppose that $m \models \bullet("A_1") \wedge \dots \wedge \bullet("A_n")$. If $\neg A$ is derivable in $O + A_1, \dots, A_n$, then A cannot be in m , otherwise $O + m$ would be inconsistent. This implies $m \not\models \bullet("A")$, that is $m \models \neg \bullet("A")$, which implies (3.1). ■

Let us now consider the combinations of reflection rules which are not characterizable axiomatically. We start with $\text{Rup}_r + \text{Rdw}_r$. For any set x of interpretations of L_M we define:

$$\begin{aligned} \bigcap x &= \{A \in L_O : \text{for all } m \in x, A \in m\} \\ \bigcup x &= \{A \in L_O : \text{there exists an } m \in x \text{ } A \in m\}. \end{aligned}$$

THEOREM 3.3

Let OM be an OM pair composed of O and M connected by $\text{Rup}_r + \text{Rdw}_r$. The set of models of $\text{TH}_{\text{OM}}(M)$ is the largest subset of interpretations of L_M that satisfies (3.2) in x :

$$x = \{m \models \Omega_M : \text{TH}(O + \bigcap x) \subseteq m\}. \quad (3.2)$$

Theorem 3.3 states that m is a model of $\text{TH}_{\text{OM}}(M)$ iff it satisfies Ω_M and contains the set of formulas which are logical consequences of the intersection of all the models of $\text{TH}_{\text{OM}}(M)$. This theorem does not provide an explicit definition of the set of models of $\text{TH}_{\text{OM}}(M)$.² Indeed, to prove that a certain interpretation m is a model of $\text{TH}_{\text{OM}}(M)$ by exploiting Theorem 3.3, one needs, first to prove that a certain set \mathfrak{M} of interpretations is contained in the set of models of $\text{TH}_{\text{OM}}(M)$ (by showing that \mathfrak{M} is a solution of (3.2)), and then to prove that $m \in \mathfrak{M}$. Let us see some examples.

EXAMPLE 3.4

If Ω_M is the empty set, then a set \mathfrak{M} which satisfies (3.2), is the set of $m \subseteq L_O$ which contains $\text{TH}(O)$. Indeed any $m \in \mathfrak{M}$ satisfies an empty set of axioms Ω_M , and $\text{TH}(O + \bigcap \mathfrak{M}) = \text{TH}(O) \subseteq m$, because $\bigcap \mathfrak{M} = \text{TH}(O)$.

EXAMPLE 3.5

If Ω_M contains the single axiom $\bullet("p")$, then a set \mathfrak{M} that satisfies (3.2) is the set of $m \subseteq L_O$ which contains $\text{TH}(O + p)$.

EXAMPLE 3.6

If Ω_M contains the single axiom $\bullet("p") \vee \bullet("q")$, then a set \mathfrak{M} that satisfies (3.2) is the set of $m \subseteq L_O$ such that $\text{TH}(O) \cup p \subseteq m$ or $\text{TH}(O) \cup q \subseteq m$.

²This is the model theoretic counterpart of the fixpoint characterization of the metatheory generated by $\text{Rup}_r + \text{Rdw}_r$.

EXAMPLE 3.7

If Ω_M is the axiom $\neg \bullet$ (“ p ”), then a set \mathfrak{M} that satisfies (3.2) is the set of subsets of L_O which contain $\text{TH}(O)$ and does not contain p .

In the previous examples, notice that the models of $\text{TH}_{\text{OM}}(M)$ might not be closed under logical consequence. To force this closure we have to relax the restriction on the application of reflection up.

PROOF. [Theorem 3.3] The proof is in two steps. In the first step we prove that the set \mathfrak{M} of models of the metatheory generated by $\text{Rup}_r + \text{Rdw}_r$ satisfies (3.2). In the second step we prove that any other solution of (3.2) is contained in \mathfrak{M} .

First step

From [2, Theorem 4.4] we have that $\text{TH}_{\text{OM}}(O)$ and $\text{TH}_{\text{OM}}(M)$ are the smallest sets of wffs which satisfy the following equations:

$$x_O = \text{TH}(O + \text{Rdw}(x_M)) \quad (3.3)$$

$$x_M = \text{TH}(M + \text{Rup}(x_O)). \quad (3.4)$$

We have therefore that

$$\text{TH}_{\text{OM}}(M) = \text{TH}(M + \text{Rup}(\text{TH}(O + \text{Rdw}(\text{TH}_{\text{OM}}(M)))).$$

Let \mathfrak{M} be the set of models of $\text{TH}_{\text{OM}}(M)$. $m \in \mathfrak{M}$ if and only if it satisfies Ω_M and satisfies all the formulas in $\text{Rup}(\text{TH}(O + \text{Rdw}(\text{TH}_{\text{OM}}(M))))$. We have therefore that $\text{TH}(O + \text{Rdw}(\text{TH}_{\text{OM}}(M))) \subseteq m$, and that

$$\mathfrak{M} = \{m \models \Omega_M : \text{TH}(O + \text{Rdw}(\text{TH}_{\text{OM}}(M))) \subseteq m\}.$$

Since \bullet (“ A ”) belongs to $\text{TH}_{\text{OM}}(M)$ if and only if for all $m \in \mathfrak{M}$, $A \in m$, we can conclude that $\text{Rdw}(\text{TH}_{\text{OM}}(M)) = \bigcap \mathfrak{M}$. We have therefore that:

$$\mathfrak{M} = \{m \models \Omega_M : \text{TH}(O + \bigcap \mathfrak{M}) \subseteq m\}.$$

Second step

Let us prove that the \mathfrak{M} is the largest set of interpretations of L_M which satisfies (3.2), namely, that if \mathfrak{N} is a solution of (3.2), then $\mathfrak{N} \subseteq \mathfrak{M}$. To do this let us consider the following sequence of theories M_0, M_1, \dots (originally defined in the proof of [2, Theorem 4.4]):

$$M_0 = M \quad (3.5)$$

$$O_0 = O \quad (3.6)$$

$$M_{i+1} = M_i + \text{Rup}(\text{TH}(O_i)) \quad (3.7)$$

$$O_{i+1} = O_i + \text{Rdw}(\text{TH}(M_i)). \quad (3.8)$$

For any $i \geq 0$, let \mathfrak{M}_i be the set of models of M_i . We prove by induction on i that for all $i \geq 0$, $\mathfrak{N} \subseteq \mathfrak{M}_i$.

Base case

$$\mathfrak{N} \subseteq \{m : m \models \Omega_M\} = \mathfrak{M}_0.$$

Step case

Since \mathfrak{N} satisfies (3.2),

$$\mathfrak{N} = \{m \models \Omega_M : \text{TH}(O + \bigcap \mathfrak{N}) \subseteq m\}.$$

The induction hypothesis $\mathfrak{N} \subseteq \mathfrak{M}_i$ implies that $\bigcap \mathfrak{M}_i \subseteq \bigcap \mathfrak{N}$, and therefore that

$$\mathfrak{N} \subseteq \{m \models \Omega_M : \text{TH}(O + \bigcap \mathfrak{M}_i) \subseteq m\}.$$

Since \mathfrak{M}_i is the set of models of M_i , we have that $\bigcap \mathfrak{M}_i = \text{Rdw}(\text{TH}(M_i))$. If we substitute $\bigcap \mathfrak{M}_i$ with $\text{Rdw}(\text{TH}(M_i))$ in the previous equation we obtain:

$$\mathfrak{N} \subseteq \{m \models \Omega_M : \text{TH}(O + \text{Rdw}(\text{TH}(M_i))) \subseteq m\}.$$

Notice that any interpretation m that satisfies Ω_M and $\text{Rup}(O + \text{Rdw}(\text{TH}(M_i)))$ is a model of M_{i+1} , and therefore, we have that

$$\mathfrak{N} \subseteq \mathfrak{M}_{i+1}.$$

Thus we have shown that for any $i \geq 0$, $\mathfrak{N} \subseteq \mathfrak{M}_i$. This implies that:

$$\mathfrak{N} \subseteq \bigcap_{i \geq 0} \mathfrak{M}_i.$$

Since $\text{TH}_{\text{OM}}(M) = \bigcup_{i \geq 0} M_i$, the set of models of M is the intersection of the models of each M_i , i.e. $\mathfrak{M} = \bigcap_{i \geq 0} \mathfrak{M}_i$. Thus we conclude that $\mathfrak{N} \subseteq \mathfrak{M}$. \blacksquare

The set of models of the metatheories generated by $\text{Rup}_r^n + \text{Rdw}_r$ and $\text{Rup}_r + \text{Rdw} + \text{Rdw}_r^n$ can be characterized similarly to $\text{Rup}_r + \text{Rdw}_r$. For a set of wff Γ , let $\neg\Gamma$ be defined as $\neg\Gamma = \{\neg A : A \in \Gamma\}$.

THEOREM 3.8

Let OM be an OM pair composed of O , M connected by $\text{Rup}_r^n + \text{Rdw}_r$. The set of models of $\text{TH}_{\text{OM}}(M)$ is the largest subset of interpretations of L_M which satisfies (3.9) in x :

$$x = \{m \models \Omega_M : \text{TH}(O + \bigcap x) \cap \neg m \neq \emptyset\}. \quad (3.9)$$

THEOREM 3.9

Let OM be an OM pair composed of O , M connected by $\text{Rup}_r + \text{Rdw} + \text{Rdw}_r^n$. The set of models of $\text{TH}_{\text{OM}}(M)$ is the largest subset of interpretations of L_M which satisfies (3.10) in x :

$$x = \left\{ m \models \Omega_M : \begin{array}{l} m \text{ is the set of theorems of an extension of } O \\ \text{If } A \notin \bigcup x, \text{ then } \neg A \in m \end{array} \right\} \quad (3.10)$$

The proofs of these two theorems are analogous to that of Theorem 3.3.

Rdw_r^n implements a weak form of negation as failure. The models of the metatheory generated by $\text{Rup}_r + \text{Rdw} + \text{Rdw}_r^n$ indeed are the extensions of O , which satisfy Ω_M and contain $\neg A$ whenever A is not provable in any extension of O .

3.2 The object theory in terms of the models of M

Our goal here is to understand the meaning of \bullet in terms of properties of the object theory, i.e. which is the property of the object theory represented by \bullet . This allows us to make precise the sense in which the object theory can be considered a model of the metatheory.

We start from the weakest combinations of reflection rules and then move to those of increasing strength.

THEOREM 3.10

Let OM be an OM pair composed of O and M connected by the set of reflection rules (RR) . Let \mathfrak{M} be the set of models of $TH_{OM}(M)$. Then:

1. If $(RR) =_O \text{Rdw}_r$, then $TH_{OM}(O) = TH(O + \bigcap \mathfrak{M})$.
2. If $(RR) =_O \text{Rup}_r$, then $TH_{OM}(O) = TH(O)$.

The underlying intuitions are the same as those described above for Theorem 3.1.

THEOREM 3.11

Let OM be an OM pair composed of O and M connected by a set of reflection rules (RR) . Let \mathfrak{M} be the set of models of $TH_{OM}(M)$. Then:

1. If $\text{Rdw}_r \leq_O (RR)$, then $\bigcap \mathfrak{M} \subseteq TH_{OM}(O)$.
2. If $\text{Rup}_r \leq_M (RR)$, then $TH_{OM}(O) \subseteq \bigcap \mathfrak{M}$.
3. If $\text{Rdw}_r \leq_O (RR)$ and $\text{Rup}_r \leq_M (RR)$, then $TH_{OM}(O) = \bigcap \mathfrak{M}$.

PROOF. [Theorem 3.11] Item (1) is a straightforward consequence of item (1) of Theorem 3.10. Item (2) is provable as follows. Suppose that $A \in TH_{OM}(O)$, then, by Rup_r , $\bullet("A") \in TH_{OM}(M)$ and therefore, $A \in \bigcap \mathfrak{M}$. Item (3) is the conjunction of (1) and (2). ■

If the object theory is generated by a set of reflection rules weaker than Rup_r or Rdw_r (in the sense of Theorem 3.11), then the set of object theorems is not completely determined by the metatheory. Vice versa, from Theorem 3.11 we can conclude that the object theory (more precisely, its set of theorems) generated by a set of reflection rules which are at least as strong as Rup_r and Rdw_r , is completely determined by the set of models (theorems) of the metatheory. Finally notice that in general, according to Theorem 3.11, an interpretation m of L_M with $m = TH_{OM}(O)$ is not necessarily a model of $TH_{OM}(M)$. Consider the following example.

EXAMPLE 3.12

Let OM be composed of the empty object theory and the metatheory with a set of axioms $\Omega_M = \{\bullet("p") \vee \bullet("q")\}$, connected by $\text{Rup}_r + \text{Rdw}$. A model of the metatheory of OM contains either p or q . In particular $m_1 = TH(O + p)$ and $m_2 = TH(O + q)$ are both models of the metatheory of OM . On the other hand the intersection of these models (which is $TH_{OM}(O)$) does not contain p , nor q . Therefore, $TH_{OM}(O)$ is not a model of $TH_{OM}(M)$.

Theorem 3.11 provides the characterization of the object theory for all the combinations of reflection rules with the exception of $\text{Rup}_r^n + \text{Rdw}$ and $\text{Rup}_r^n + \text{Rdw}$. For these two special cases we have the following result.

THEOREM 3.13

Let OM be an OM pair composed of O and M connected by the set of reflection rules (RR) . If \mathfrak{M} is the set of models of $TH_{OM}(M)$, then:

1. If (RR) is $\text{Rup}_r^n + \text{Rdw}$, then $\text{TH}_{\text{OM}}(O) = \text{TH}(O + \bigcap \mathfrak{M})$;
2. If (RR) is $\text{Rup}^n + \text{Rdw}$, then $\text{TH}_{\text{OM}}(O) = \bigcap_{m \in \mathfrak{M}} \text{TH}(O + m)$.

PROOF. [Theorem 3.13] Item (1). From [2, Theorem 4.3] we have

$$\text{TH}_{\text{OM}}(O) = \text{TH}(O + \text{Rdw}(\text{TH}_{\text{OM}}(M))). \quad (3.11)$$

Since $\text{Rdw}(\text{TH}_{\text{OM}}(M)) = \bigcap \mathfrak{M}$, (3.11) implies that $\text{TH}_{\text{OM}}(O) = \text{TH}(O + \bigcap \mathfrak{M})$. To prove item (2), by [2, Theorem 4.3], we have

$$\text{TH}_{\text{OM}}(O) = \text{TH}(O + \{\bigvee_{i=1}^n A_i : \bigvee_{i=1}^n \bullet("A_i") \in \text{TH}_{\text{OM}}(M)\}). \quad (3.12)$$

To prove that $\text{TH}_{\text{OM}}(O) \subseteq \bigcap_{m \in \mathfrak{M}} \text{TH}(O + m)$, we show that for any extra axiom $A_1 \vee \dots \vee A_n$ added to O to obtain $\text{TH}_{\text{OM}}(O)$, we have that $A_1 \vee \dots \vee A_n \in \bigcap_{m \in \mathfrak{M}} \text{TH}(O + m)$. For each $m \in \mathfrak{M}$, if $\bullet("A_1") \vee \dots \vee \bullet("A_n") \in \text{TH}_{\text{OM}}(M)$ then $m \models \bullet("A_1") \vee \dots \vee \bullet("A_n")$, which implies that there is an $1 \leq i \leq n$ such that $A_i \in m$, and $A_1 \vee \dots \vee A_n \in \text{TH}(O + m)$. Therefore, for each $m \in \mathfrak{M}$, $\text{TH}_{\text{OM}}(O) \subseteq \text{TH}(O + m)$; and therefore, $\text{TH}_{\text{OM}}(O) \subseteq \bigcap_{m \in \mathfrak{M}} \text{TH}(O + m)$.

Vice versa, let us prove $\bigcap_{m \in \mathfrak{M}} \text{TH}(O + m) \subseteq \text{TH}_{\text{OM}}(O)$. Let $A \in \bigcap_{m \in \mathfrak{M}} \text{TH}(O + m)$, and for each $m \in \mathfrak{M}$ let Π_m be a proof of A in $O + m$. Let M' be the following theory:

$$M + \text{TH}_{\text{OM}}(M) + \left\{ \begin{array}{l} \neg(\wedge \bullet(" \Gamma_m ")) : m \in \mathfrak{M}, \text{ and } \Gamma_m \subseteq m \text{ is the finite (possibly} \\ \text{empty) set of formulas occurring as axioms in} \\ \Pi_m \end{array} \right\}$$

where $\wedge \bullet(" \Gamma_m ")$ is the conjunction of the formulas in $\bullet(" \Gamma_m ")$. Let us prove that M' is unsatisfiable. By contradiction let m be a model of M' . m is a model of $\text{TH}_{\text{OM}}(M)$ and therefore, $\Gamma_m \subseteq m$. This implies that $m \models (\wedge \bullet(" \Gamma_m "))$ and therefore that $m \not\models M'$. M' unsatisfiable implies that M' is finitely unsatisfiable, i.e. there is a finite sequence of sets $\Gamma_{m_1}, \dots, \Gamma_{m_n}$ such $\text{TH}_{\text{OM}}(M) + \{\neg(\wedge \bullet(" \Gamma_{m_1} ")), \dots, \neg(\wedge \bullet(" \Gamma_{m_n} "))\}$ is unsatisfiable. This implies that

$$(\wedge \bullet(" \Gamma_{m_1} ")) \vee \dots \vee (\wedge \bullet(" \Gamma_{m_n} ")) \in \text{TH}_{\text{OM}}(M)$$

and therefore that for each $\gamma_1, \dots, \gamma_n$, with $\gamma_i \in \Gamma_{m_i}$,

$$\bullet(" \gamma_1 ") \vee \dots \vee \bullet(" \gamma_n ") \in \text{TH}_{\text{OM}}(M).$$

We have therefore that

$$\gamma_1 \vee \dots \vee \gamma_n \in \text{TH}_{\text{OM}}(O).$$

A possible proof of A in $\text{TH}_{\text{OM}}(O)$ is therefore the following:

$$\frac{\frac{\frac{\vdots}{\bigwedge_{\gamma_i \in \Gamma_i} (\gamma_1 \vee \dots \vee \gamma_n)}{(\wedge \Gamma_1) \vee \dots \vee (\wedge \Gamma_n)} \quad \frac{[\wedge \Gamma_1]}{\Pi_1} \quad \frac{[\wedge \Gamma_n]}{\Pi_n}}{A} \quad \dots \quad A}{A} \text{VE}$$

We can therefore conclude that $A \in \text{TH}_{\text{OM}}(O)$. ■

TABLE 1. Metamodels and object theorems generated by a set of reflection rules

(RR)	Models of $\text{TH}_{\text{OM}}(M)$	$\text{TH}_{\text{OM}}(O)$
\emptyset	Any set of L_O -wffs	$\text{TH}(O)$
Rdw_r	Any set of L_O -wffs	$\text{TH}(O + \bigcap \mathfrak{M})$
Rup_r	$\text{TH}(O) \subseteq m$	$\text{TH}(O)$
Rdw	Any set of L_O -wffs	$\text{TH}(O + \bigcap \mathfrak{M})$
Rup	$\text{TH}(O) \subseteq m$	$\text{TH}(O)$
$\text{Rup}_r + \text{Rdw}_r$	Fixpoint equation (3.2)	$\bigcap \mathfrak{M}$
$\text{Rup}_r^n + \text{Rdw}$	m is consistent with $\text{TH}(O)$	$\text{TH}(O + \bigcap \mathfrak{M})$
$\text{Rup}_r + \text{Rdw}$	Any extension of O	$\bigcap \mathfrak{M}$
$\text{Rup}_r + \text{Rdw} + \text{Rup}_r^n$	Any consistent extension of O	$\bigcap \mathfrak{M}$
$\text{Rup}_r + \text{Rdw} + \text{Rdw}_r^n$	fixpoint equation (3.10)	$\bigcap \mathfrak{M}$
$\text{Rup}_r + \text{Rdw} + \text{Rup}^n$	Any consistent extension of O	$\bigcap \mathfrak{M}$
$\text{Rup}_r + \text{Rdw} + \text{Rdw}^n$	Any maximal extension of O	$\bigcap \mathfrak{M}$
$\text{Rup}^n + \text{Rdw}$	m is consistent with $\text{TH}(O)$	$\bigcap_{m \in \mathfrak{M}} \text{TH}(O + m)$
$\text{Rup} + \text{Rdw}$	Any maximal extension of O	$\bigcap \mathfrak{M}$
$\text{Rup} + \text{Rdw} + \text{Rup}_r^n$	Any maximal consistent extension of O	$\bigcap \mathfrak{M}$

Notice that, $\text{Rup}_r^n + \text{Rdw}$ and $\text{Rup}^n + \text{Rdw}$ generate the same metatheory but they do not generate the same object theory. Consider the following example.

EXAMPLE 3.14

Let OM be the OM pair composed of the empty object theory O and the metatheory M with axiom $\bullet("p") \vee \bullet("q")$, connected by the reflection rules $\text{Rup}^n + \text{Rdw}$. Let OM_r be the OM pair obtained from OM by adding the restriction to Rup^n . Since they have the same set of metatheorems, OM_r and OM have the same set of models \mathfrak{M} defined as follows:

$$\mathfrak{M} = \{m : m \text{ is consistent and } p \in m \text{ or } q \in m\}. \quad (3.13)$$

The set of object theorems of OM_r differs from that of OM . Indeed, according to Theorem 3.13, we have the following. The set of object level theorems of OM_r is equal to the theorems provable from O extended with the intersection of all the elements of \mathfrak{M} , which is the empty set. Therefore, $\text{TH}_{\text{OM}_r}(O) = \text{TH}(O)$ is the set of the propositional tautologies. On the other hand the set of object level theorems of OM is equal to the intersection of the theorems provable from O extended with m for some $m \in \mathfrak{M}$. This means, for instance, that, since either p or q belongs to any m , $p \vee q$ is provable from $O + m$ for all m . This implies that $p \vee q \in \text{TH}_{\text{OM}}(O)$.

Table 1 summarizes the results proved in this section.

4 A strict ordering characterization

In Part I [2] reflection rules have been partially ordered; however, we still do not have explicit results stating that if $(\text{RR})_i \leq_X (\text{RR})_2$, it is not the case that $(\text{RR})_i =_X (\text{RR})_2$. In this section we show that all the combinations of reflection rules described in this paper generate

different object-meta relations, and therefore the partial order described in the previous paper can be transformed in a *strict partial order*.

Let us recall the definition of partial orders between combinations of reflection rules given in the previous paper.

DEFINITION 4.1

For each two OM pairs OM_1 and OM_2 ,

1. $OM_1 \leq_O OM_2$ if $TH_{OM_1}(O) \subseteq TH_{OM_2}(O)$;
2. $OM_1 \leq_M OM_2$ if $TH_{OM_1}(M) \subseteq TH_{OM_2}(M)$;
3. $OM_1 \leq_D OM_2$ if $\vdash_{OM_1} \subseteq \vdash_{OM_2}$;
4. $OM_1 =_X OM_2$ if $OM_1 \leq_X OM_2$ and $OM_2 \leq_X OM_1$, where $X \in \{O, M, D\}$;
5. $OM_1 <_X OM_2$ if $OM_1 \leq_X OM_2$ and not $OM_2 =_X OM_1$, where $X \in \{O, M, D\}$.

Orders on sets of reflection rules can be defined on the basis on the orders of OM pairs.

DEFINITION 4.2

Let $OM(RR)$ be an OM pair composed of an object theory and a metatheory connected by the set of reflection rules (RR) . For any combination of reflection rules $(RR)_1, (RR)_2$:

1. $(RR)_1 \leq_O (RR)_2$ if, for any $OM(RR)_1$ and $OM(RR)_2$, with the same object theory and metatheory, $OM(RR)_1 \leq_O OM(RR)_2$;
2. $(RR)_1 \leq_M (RR)_2$ if, for any $OM(RR)_1$ and $OM(RR)_2$ with the same object theory and metatheory, $OM(RR)_1 \leq_M OM(RR)_2$;
3. $(RR)_1 \leq_D (RR)_2$ if, for any $OM(RR)_1$ and $OM(RR)_2$ with the same object theory and metatheory, $OM(RR)_1 \leq_D OM(RR)_2$;
4. For any $X \in \{O, M, D\}$, $(RR)_1 =_X (RR)_2$, if both $(RR)_1 \leq_X (RR)_2$ and $(RR)_2 \leq_X (RR)_1$.
5. For any $X \in \{O, M, D\}$, $(RR)_1 <_X (RR)_2$, if $(RR)_1 \leq_X (RR)_2$ and not $(RR)_2 =_X (RR)_1$.

The orders defined above are not completely independent of one another. For instance, if two combinations of reflection rules are equivalent (i.e. $(RR)_1 =_D (RR)_2$) then they generate the same object theory and metatheory (i.e. $(RR)_1 =_O (RR)_2$ and $(RR)_1 =_M (RR)_2$). The relations among the three orders are summarized in Figure 1. In this figure, an arrow connects X to Y , if and only if $(RR)_1 X (RR)_2$ implies $(RR)_1 Y (RR)_2$. The correctness of this latter graph is straightforward.

THEOREM 4.3 (Correctness of the graph in Figure 2)

If the graph of Figure 2 contains an arc labelled with X from $(RR)_1$ to $(RR)_2$ then $(RR)_1 X (RR)_2$.

To prove the previous theorem we need the following lemma:

LEMMA 4.4 (Effective assumptions)

Let OM be an OM pair with the set (RR) of reflection rules. Let Γ_O, A_O and Γ_M, A_M be a set of wffs of L_O and L_M respectively. Then:

1. if the reflection up rules in (RR) are restricted, then $\Gamma_O, \Gamma_M \vdash_{OM} A_M$ if and only if $\Gamma_M \vdash_{OM} A_M$;

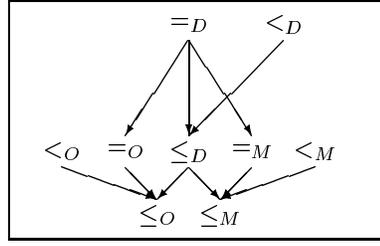


FIGURE 1. Relations among orders

2. if the reflection down rules in (RR) are restricted, then $\Gamma_O, \Gamma_M \vdash_{\text{OM}} A_O$ if and only if $\Gamma_O \vdash_{\text{OM}} A_O$.

To prove (1), it is enough to observe that only way to ‘switch’ from object to meta in order to infer a metaformula from object assumptions, is by applying restricted reflection rules. However, these rules are applicable only if all the assumptions in the object theory are discharged. A similar argument can be used to prove (2).

PROOF. [Theorem 4.3] The graph in Figure 2 is obtained from the analogous graph in [2] by refining the relations between reflection rules. Since $<_O$ (resp. $<_M$ and $<_D$) are defined as the conjunction of \leq_D and \neq_O (resp. the conjunction of \leq_D and \neq_M and the conjunction of \leq_D and \neq_D), the additional information encoded in the graph in Figure 2, with respect to the previous version, consists of statements asserting that two combinations of reflection rules differ either in the object and/or metatheory they generate, or in their derivability relation.

Following Definition 4.2, we prove $(\text{RR})_1 \neq_O (\text{RR})_2$ (resp. $(\text{RR})_1 \neq_M (\text{RR})_2$), by showing that there are two OM pairs OM_1 and OM_2 composed of the same object and metatheories and connected by $(\text{RR})_1$ and $(\text{RR})_2$ respectively, which generate a different object (resp. meta) theory. In practice, we look for an object (resp. meta) wff A not provable in OM_1 and provable in OM_2 . $(\text{RR})_1 \neq_D (\text{RR})_2$ is automatically proved as a consequence of $(\text{RR})_1 \neq_O (\text{RR})_2$ or $(\text{RR})_1 \neq_M (\text{RR})_2$. We prove \neq_D without proving $(\text{RR})_1 \neq_O (\text{RR})_2$ or $(\text{RR})_1 \neq_M (\text{RR})_2$ by providing two OM pairs OM_1 and OM_2 and a set of object and meta wffs Γ, A such that A is not derivable from Γ in OM_1 , and A is derivable from Γ in OM_2 .

We proceed by considering the arcs of the graph in Figure 2 from the bottom up. We say that a (meta or object) theory is an empty (meta or object) theory meaning that it has no axioms.

1. $\emptyset <_M \text{Rup}_r$. (We prove $\emptyset \neq_M \text{Rup}_r$.) Consider the OM pairs OM_1 and OM_2 composed of the empty object theory and the empty metatheory connected by no reflection rules and by Rup_r respectively. No wff of the form $\bullet(“A”)$ is a theorem of the metatheory of OM_1 , while, for instance, $\bullet(“\neg \perp”)$ is a metatheorem of OM_2 .
2. $\emptyset <_O \text{Rdw}_r$. (We prove $\emptyset \neq_O \text{Rdw}_r$.) Consider the OM pairs OM_1 and OM_2 composed of the empty object theory and the metatheory with the only axiom $\bullet(“\perp”)$ connected by no reflection rules and by Rdw_r respectively. Since the object theory of OM_1 is equal to O , \perp is not provable, Instead \perp is provable (by Rdw_r) in the object theory of OM_2 .
3. $\text{Rup}_r <_D \text{Rup}$. (We prove $\text{Rup}_r \neq_D \text{Rup}$.) Consider two OM pairs OM_1 and OM_2 composed of the empty object theory O and the empty metatheory M , connected by Rup_r and Rup respectively. The metatheory of OM_1 is equal to M and therefore $\bullet(“\perp”)$ is not

provable in it. Furthermore, by Lemma 4.4, $\bullet(" \perp ")$ is not derivable in the metatheory from the assumption \perp of the object theory. In OM_2 instead, $\bullet(" \perp ")$ is derivable (by Rup) from the assumption \perp in the object theory.

4. $\text{Rdw}_r <_D \text{Rdw}$. Analogous to the previous case.
5. $\text{Rup}_r <_M \text{Rup}_r + \text{Rdw}_r$. (We prove $\text{Rup}_r \neq_M \text{Rup}_r + \text{Rdw}_r$.) Consider two OM pairs OM_1 and OM_2 composed of the empty object theory O and the metatheory M with the axiom $\bullet(" \perp ")$ connected by Rup_r and $\text{Rup}_r + \text{Rdw}_r$ respectively. By Theorem 3.1, the interpretation $m = \text{TH}(O) \cup \{\perp\}$ is a model of the metatheory of OM_1 . As O is consistent, there is an object wff A distinct from \perp such that $m \not\models \bullet("A")$ and therefore $\bullet("A")$ is not a theorem of the metatheory of OM_1 . In OM_2 , instead, each formula $\bullet("A")$ is provable as follows:

$$\frac{\frac{\bullet(" \perp ") \quad \text{Rdw}_r}{\frac{\perp}{A} \perp}}{\bullet("A")} \text{Rup}_r$$

6. $\text{Rup}_r <_O \text{Rup}_r + \text{Rdw}_r$. (We prove $\text{Rup}_r \neq_O \text{Rup}_r + \text{Rdw}_r$.) In the above proof we can see that \perp is provable in the object theory of OM_2 while it is not provable in the object theory of OM_1 , because there is no rule going from the metatheory to the object theory.
7. $\text{Rdw}_r <_O \text{Rup}_r + \text{Rdw}_r$. (We prove $\text{Rdw}_r \neq_O \text{Rup}_r + \text{Rdw}_r$.) Consider two OM pairs OM_1 and OM_2 composed of the empty object theory O and the metatheory M with the axiom $\neg \bullet(" \neg \perp ")$ connected by Rdw_r and $\text{Rup}_r + \text{Rdw}_r$ respectively. From [2, Theorem 4.1] we have that the set of object theorems of OM_1 is $\text{TH}(O + \text{Rdw}(\text{TH}(M)))$. Since M does not prove any theorem of the form $\bullet("A")$, the set of object theorems of OM_1 is $\text{TH}(O)$. In particular \perp is not provable in O and therefore it is not provable in OM_1 . In OM_2 instead \perp is a theorem of the object theory. Indeed:

$$\frac{\frac{\frac{\perp}{\neg \perp} \supset \text{I}}{\bullet(" \neg \perp ") \quad \text{Rup}_r} \neg \bullet(" \neg \perp ") \quad \supset \text{E}}{\frac{\frac{\perp}{\bullet(" \perp ") \quad \perp}}{\perp} \text{Rdw}_r} \supset \text{E}$$

8. $\text{Rdw}_r <_M \text{Rup}_r + \text{Rdw}_r$. (We prove $\text{Rdw}_r \neq_M \text{Rup}_r + \text{Rdw}_r$.) Consider the OM pairs OM_1 , and OM_2 composed by the empty object theory and metatheory connected by Rup_r and $\text{Rup}_r + \text{Rdw}_r$ respectively. $\bullet(" \top ")$ is not provable in OM_1 as the empty set is a model of the metatheory of OM_1 , while $\bullet(" \top ")$ is provable in OM_2 .
9. $\text{Rup} <_D \text{Rup} + \text{Rdw}$. (We prove $\text{Rup} \neq_D \text{Rup} + \text{Rdw}$.) Consider two OM pairs OM_1 and OM_2 composed of the empty object theory O and the empty metatheory M , connected by Rup and $\text{Rup} + \text{Rdw}$ respectively. By [2, Theorem 4.1] we have that the set of theorems of the object theory of OM_1 is $\text{TH}(O)$ and therefore \perp is not provable in the object theory of OM_1 . Therefore, by Lemma 4.4, in OM_1 \perp is not derivable in the object theory from $\bullet(" \perp ")$. In OM_2 , instead, \perp can be derived in the object theory from $\bullet(" \perp ")$ in the metatheory (via a single application of Rdw).
10. $\text{Rup}_r + \text{Rdw}_r <_M \text{Rup}_r + \text{Rdw}$. (We prove $\text{Rup}_r + \text{Rdw}_r \neq_M \text{Rup}_r + \text{Rdw}$.) Consider two OM pairs OM_1 and OM_2 composed of the empty object theory O and the empty

metatheory M , connected by $\text{Rup}_r + \text{Rdw}_r$ and $\text{Rup}_r + \text{Rdw}$ respectively. By Theorem 3.3 the set of models of the metatheory of OM_1 is $\{m : \text{TH}(O) \subseteq m\}$. Consider the model of the metatheory of OM_1 $m = \text{TH}(O) \cup \{\perp\}$. Clearly $m \not\models \bullet(\text{"}\perp\text{"}) \supset \bullet(\text{"}p\text{"})$, and therefore $\bullet(\text{"}\perp\text{"}) \supset \bullet(\text{"}p\text{"})$ is not a theorem of OM_1 . However, $\bullet(\text{"}\perp\text{"}) \supset \bullet(\text{"}p\text{"})$ is a theorem of OM_2 .

11. $\text{Rup}_r + \text{Rdw}_r <_O \text{Rup}_r + \text{Rdw}$. (We prove $\text{Rup}_r + \text{Rdw}_r \neq_O \text{Rup}_r + \text{Rdw}$.) Consider two OM pairs OM_1 and OM_2 composed of the empty object theory O and the metatheory M with axiom $\bullet(\text{"}p\text{"}) \vee \bullet(\text{"}q\text{"})$, connected by $\text{Rup}_r + \text{Rdw}_r$ and $\text{Rup}_r + \text{Rdw}$, respectively. By Theorem 3.3 the set of models of the metatheory of OM_1 is:

$$\{m : \text{TH}(O) \subseteq m \text{ and } p \in m \text{ or } q \in m\}. \quad (4.1)$$

By Theorem 3.11 the set of theorems of the object theory of OM_1 is the intersection of m for each m in (4.1), which is $\text{TH}(O)$. The formula $p \vee q$ is not provable in $\text{TH}(O)$ as O , being empty, proves only classical tautologies. This implies that $p \vee q$ is not a theorem of the object theory of OM_1 . In OM_2 , instead, $p \vee q$ is provable as follows:

$$\frac{\frac{\frac{\bullet(\text{"}p\text{"})}{p} \text{Rdw} \quad \frac{\bullet(\text{"}q\text{"})}{q} \text{Rdw}}{p \vee q} \text{VI} \quad \frac{\bullet(\text{"}p\text{"}) \vee \bullet(\text{"}q\text{"})}{\bullet(\text{"}p \vee q\text{"})} \text{Rup}_r}{\frac{\bullet(\text{"}p \vee q\text{"})}{p \vee q} \text{Rdw}} \frac{\frac{\bullet(\text{"}p \vee q\text{"})}{\bullet(\text{"}p \vee q\text{"})} \text{Rup}_r}{\bullet(\text{"}p \vee q\text{"})} \text{VE} \quad \text{VE}$$

12. $\text{Rdw} <_M \text{Rup}_r^n + \text{Rdw}$. (We prove $\text{Rdw} \neq_M \text{Rup}_r^n + \text{Rdw}$.) Consider two OM pairs OM_1 and OM_2 composed of the empty object theory O and the empty metatheory M . By [2, Theorem 4.1], the set of theorems of the metatheory of OM_1 is $\text{TH}(M)$ (i.e. the set of tautologies of L_M), and by [2, Theorem 4.1], the set of metatheorems of OM_2 is $\text{TH}(M + (\text{Cons}))$ which contains wffs which are not tautologies of L_M (e.g. $\bullet(\text{"}p\text{"}) \wedge \bullet(\text{"}q\text{"}) \supset \neg \bullet(\text{"}\neg p \vee \neg q\text{"})$).
13. $\text{Rdw} <_O \text{Rup}_r^n + \text{Rdw}$. (We prove $\text{Rdw} \neq_O \text{Rup}_r^n + \text{Rdw}$.) Consider two OM pairs OM_1 and OM_2 composed of the empty object theory O and the metatheory M with the axiom $\bullet(\text{"}\perp\text{"}) \vee \bullet(\text{"}p\text{"})$. By Theorem 3.1 the set of models of the metatheory of OM_1 is the set of models of M , i.e. the set of models m such that $\perp \in m$ or $p \in m$. By Theorem 3.11 the set of object theorems of OM_1 is the set of theorems of O extended with the intersection of the models of M . Since this intersection is empty, the object theory of OM_1 is equal to O and, therefore, it does not prove p . In the object theory of OM_2 , instead, p is provable as follows:

$$\frac{\frac{\frac{\frac{\perp}{\neg \perp} \supset \text{I}}{\bullet(\text{"}\perp\text{"}) \rightarrow \bullet(\text{"}\perp\text{"})} \text{Rup}_r^n}{\bullet(\text{"}\perp\text{"}) \vee \bullet(\text{"}p\text{"})} \supset \text{E} \quad \frac{\frac{\perp}{\bullet(\text{"}p\text{"})} \perp}{\bullet(\text{"}p\text{"})} \text{VE}}{\frac{\bullet(\text{"}p\text{"})}{p} \text{Rdw}} \text{VE}$$

14. $\text{Rup}_r + \text{Rdw} <_M \text{Rup}_r + \text{Rdw} + \text{Rdw}_r^n$. (We prove $\text{Rup}_r + \text{Rdw} \neq_M \text{Rup}_r + \text{Rdw} + \text{Rdw}_r^n$.) Consider two OM pairs OM_1 and OM_2 composed of the empty object

theory O and the metatheory M , with the axiom $\neg \bullet ("p")$. Since p is not provable from the empty object theory, the interpretation $m = \text{TH}(O)$ is a model of the metatheory (according to Theorem 3.2). Since also $\neg p$ is not provable in the empty object theory (i.e. $\neg p \notin \text{TH}(O)$), $m \not\models \bullet (" \neg p ")$ and therefore $\bullet (" \neg p ")$ is not provable in the metatheory of OM_1 . In OM_2 instead $\bullet (" \neg p ")$ is provable. A proof of $\bullet (" \neg p ")$ is the following:

$$\frac{\frac{\neg \bullet ("p")}{\neg p} \text{Rdw}_r^n}{\bullet (" \neg p ") \text{Rup}_r}$$

15. $\text{Rup}_r + \text{Rdw} <_O \text{Rup}_r + \text{Rdw} + \text{Rdw}_r^n$. (We prove $\text{Rup}_r + \text{Rdw} \neq_O \text{Rup}_r + \text{Rdw} + \text{Rdw}_r^n$.) Consider the two OM pairs of the previous case. The object theory of OM_1 does not prove $\neg p$. Indeed, by Theorem 3.1, the object theory of OM_1 is the intersection of the models of the metatheory of OM_1 . The model m considered in the previous case is such that $\neg p \notin m$. Therefore $\neg p$ is not provable in the object theory of OM_1 . In OM_2 instead $\neg p$ is provable (by applying Rdw_r^n to $\neg \bullet ("p")$).
16. $\text{Rup}_r + \text{Rdw} <_M \text{Rup}_r + \text{Rdw} + \text{Rup}_r^n$. (We prove that $\text{Rup}_r + \text{Rdw} \neq_M \text{Rup}_r + \text{Rdw} + \text{Rup}_r^n$.) Consider two OM pairs OM_1 and OM_2 composed of the empty object theory O and the empty metatheory M . The metatheory of OM_1 does not prove $\neg \bullet (" \perp ")$. Indeed the model m which interprets \bullet in the set of wffs L_O , is a model of $\text{TH}_{\text{OM}_1}(M)$. The metatheory of OM_2 instead proves $\neg \bullet (" \perp ")$. A proof is the following.

$$\frac{\frac{\perp}{\neg \perp} \supset \text{I}}{\neg \bullet (" \perp ") \text{Rup}_r^n} \quad (4.2)$$

17. $\text{Rup}_r^n + \text{Rdw} <_O \text{Rup}_r^n + \text{Rdw}$. See Example 3.14.
18. $\text{Rup}_r + \text{Rdw} + \text{Rdw}_r^n <_M \text{Rup}_r + \text{Rdw} + \text{Rdw}_r^n$. (We prove $\text{Rup}_r + \text{Rdw} + \text{Rdw}_r^n \neq_M \text{Rup}_r + \text{Rdw} + \text{Rdw}_r^n$.) Consider two OM pairs OM_1 and OM_2 composed of the empty object theory O and the empty metatheory M , connected by the reflection rules $\text{Rup}_r + \text{Rdw} + \text{Rdw}_r^n$ and $\text{Rup}_r + \text{Rdw} + \text{Rdw}_r^n$ respectively. By Theorem 3.9, the set of models of the metatheory of OM_1 is equal to the set of interpretations which are the sets of theorems of an extension of O . The interpretation $m = \text{TH}(O)$ is, therefore, a model of the metatheory of OM_1 . Since neither p nor $\neg p$ is provable in O , then neither $m \models \bullet ("p")$ nor $m \models \bullet (" \neg p ")$ and therefore $m \not\models \bullet ("p") \vee \bullet (" \neg p ")$. This implies that $\bullet ("p") \vee \bullet (" \neg p ")$ is not valid, and therefore not provable, in the metatheory of OM_1 . In the metatheory of OM_2 , instead, $\bullet ("p") \vee \bullet (" \neg p ")$ is provable as follows:

$$\frac{\frac{\bullet ("p")}{\bullet ("p") \vee \bullet (" \neg p ") \text{VI}} \quad \frac{\frac{\frac{\neg \bullet ("p")}{\neg p} \text{Rdw}_r^n}{\bullet (" \neg p ") \text{Rup}_r} \quad \frac{\bullet (" \neg p ") \vee \bullet (" \neg p ") \text{VI}}{\bullet ("p") \vee \bullet (" \neg p ") \text{VE}}}{\bullet ("p") \vee \bullet (" \neg p ") \text{VE}} \quad \frac{\bullet ("p") \vee \bullet (" \neg p ") \text{VI}}{\bullet ("p") \vee \bullet (" \neg p ") \text{VE}}}{\bullet ("p") \vee \bullet (" \neg p ") \text{VE}}$$

19. $\text{Rup}_r + \text{Rdw} + \text{Rdw}_r^n <_O \text{Rup}_r + \text{Rdw} + \text{Rdw}_r^n$. (We prove $\text{Rup}_r + \text{Rdw} + \text{Rdw}_r^n \neq_O \text{Rup}_r + \text{Rdw} + \text{Rdw}_r^n$.) Consider two OM pairs OM_1 and OM_2 composed of the empty object theory O , and the metatheory M with the axiom $\bullet ("p") \vee \bullet (" \neg p ") \supset \bullet ("q")$, connected by the reflection rules $\text{Rup}_r + \text{Rdw} + \text{Rdw}_r^n$ and $\text{Rup}_r + \text{Rdw} + \text{Rdw}_r^n$

respectively. By Theorem 3.9 it can be proved that the set of models of the metatheory of OM_1 is equal to the set of interpretations m which are the set of theorems of an extension of O , and $q \in m$ if $p \in m$ or $\neg p \in m$. Since neither p nor $\neg p$ is provable in O , then the interpretation $m = TH(O)$ is a model of the metatheory of OM_1 . Theorem 3.11 implies that the set of theorems of the object theory of OM_1 is equal to the intersection of each model of the metatheory. Since q does not belong to a model m , then p is not provable in the object theory of OM_1 . In OM_2 , instead, q is provable as $\bullet("p") \vee \bullet("\neg p")$ is a theorem of the metatheory of OM_2 (see the previous example) and q can be derived by $\supset E$ and Rdw .

20. $Rup_r + Rdw + Rup_r^n <_D Rup_r + Rdw + Rup^n$. (We prove $Rup_r + Rdw + Rup_r^n \neq_D Rup_r + Rdw + Rup^n$.) Consider two OM pairs OM_1 and OM_2 composed of the empty object theory O , and the empty metatheory M , connected by the reflection rules $Rup_r + Rdw + Rup_r^n$ and $Rup_r + Rdw + Rup^n$ respectively. By Theorem 3.2 the metatheory of OM_1 has at least a model (e.g. $m = TH(O)$), and therefore it is consistent. By Lemma 4.4, the non-derivability of \perp in the metatheory implies the non-derivability of \perp in the metatheory starting from \perp in the object theory. This derivation is instead possible in OM_2 :

$$\frac{\frac{\frac{\perp}{\neg\neg\perp} \perp}{\neg\bullet("\neg\perp")} Rup^n}{\perp} \quad \frac{\frac{\frac{\perp}{\neg\perp} \supset I}{\neg\perp} \supset I}{\bullet("\neg\perp")} Rup_r}{\perp} \supset E$$

21. $Rup^n + Rdw <_M Rup_r + Rdw + Rup^n$. (We prove $Rup^n + Rdw \neq_M Rup_r + Rdw + Rup^n$.) Consider two OM pairs OM_1 and OM_2 composed of the empty object theory O , and the empty metatheory M , connected by the reflection rules $Rup^n + Rdw$ and $Rup_r + Rdw + Rup^n$ respectively. By Theorem 3.2, the interpretation $m = \emptyset$ is a model of the metatheory of OM_1 . This model does not satisfy $\bullet("\neg\perp")$. This implies that $\bullet("\neg\perp")$ is not a theorem of the metatheory of OM_1 . In OM_2 instead, $\bullet("\neg\perp")$ is provable (by Rup_r).
22. $Rup^n + Rdw <_O Rup_r + Rdw + Rup^n$. (We prove $Rup^n + Rdw \neq_O Rup_r + Rdw + Rup^n$.) Consider two OM pairs OM_1 and OM_2 composed of the empty object theory O , and the metatheory M with axiom $\bullet("\neg\perp") \supset \bullet("\perp")$, connected by the reflection rules $Rup^n + Rdw$ and $Rup_r + Rdw + Rup^n$ respectively. By Theorem 3.13 we have that the set of theorems of the object theory of OM_1 is equal to the intersection of the sets $TH(O + m)$ for each model m of M . According to Theorem 3.2 the interpretation $m = \emptyset$, is a model of the metatheory of OM_1 . This implies that the set of theorems of the object theory of OM_1 is $TH(O)$. Note that O is empty and therefore consistent. The object theory of OM_2 , instead, is inconsistent, as \perp is provable in the object theory of OM_2 .
23. $Rup_r + Rdw + Rdw^n <_D Rup_r + Rdw$. (We prove $Rup_r + Rdw + Rdw^n \neq_D Rup_r + Rdw$.) The reasoning is similar to that done for the case $Rup_r <_D Rup$.
24. $Rup_r + Rdw + Rup^n <_M Rup_r + Rdw + Rup_r^n$. (We prove that $Rup_r + Rdw + Rup^n \neq_M Rup_r + Rdw + Rup_r^n$.) Consider two OM pairs OM_1 and OM_2 composed of the empty object theory O , and the empty metatheory M , connected by the reflection rules $Rup_r + Rdw + Rup^n$ and $Rup_r + Rdw + Rup_r^n$ respectively. According to Theorem 3.2 the interpretation $m = TH(O)$ is a model of the metatheory of OM_1 . Note that, for any proposition p , neither p nor $\neg p$ is provable in O . From this we can infer that $m \not\models \bullet("p")$

and $m \not\models \bullet(\neg p)$, i.e. $m \not\models \bullet(p) \vee \bullet(\neg p)$. This implies that $\bullet(p) \vee \bullet(\neg p)$ is not provable in the metatheory of OM_1 . The same formula, however, is provable in the metatheory of OM_2 (see proof above).

25. $Rup + Rdw <_M Rup + Rdw + Rup^n$. (We prove that $Rup + Rdw \not\equiv_M Rup + Rdw + Rup^n$.) Consider two OM pairs OM_1 and OM_2 composed of the empty object theory O , and the empty metatheory M , connected by the reflection rules $Rup + Rdw$ and $Rup + Rdw + Rup^n$, respectively. According to Theorem 3.2 the interpretation $mm = L_O$ is a model of the metatheory of OM_1 . Since this model does not satisfy $\neg \bullet(\perp)$, $\neg \bullet(\perp)$ is not provable in the metatheory of OM_1 . As showed before $\neg \bullet(\perp)$ is a theorem of the metatheory of OM_2 .
26. $Rdw <_D Rup_r + Rdw$. Note that $Rdw =_O Rdw_r$ and $Rdw_r <_O Rup_r + Rdw_r$, implies that $Rdw <_O Rup_r + Rdw_r$. From the fact that $Rup_r + Rdw_r <_O Rup_r + Rdw$, we have that $Rup <_O Rup_r + Rdw$. Since $(RR)_1 \leq_D (RR)_2$ and $(RR)_1 <_O (RR)_2$ implies $(RR)_1 <_D (RR)_2$, we conclude that $Rup <_D Rup_r + Rdw$.
27. $Rup_r^n + Rdw <_O Rup_r + Rdw + Rup_r^n$. Suppose, by contradiction, that $Rup_r^n + Rdw =_O Rup_r + Rdw + Rup_r^n$. From the fact that $Rup_r + Rdw + Rup_r^n =_O Rup_r + Rdw + Rup^n$, we have that $Rup_r^n + Rdw =_O Rup_r + Rdw + Rup^n$. However, from the fact that $Rup_r^n + Rdw <_O Rup^n + Rdw$, and from the fact that $Rup^n + Rdw <_O Rup_r + Rdw + Rup^n$, we have that $Rup_r^n + Rdw <_O Rup_r + Rdw + Rup^n$ and, therefore, that $Rup_r^n + Rdw \neq_O Rup_r + Rdw + Rup^n$. This contradicts what we have proved before. We conclude therefore that $Rup_r^n + Rdw \neq_O Rup_r + Rdw + Rup_r^n$.
28. $Rup_r^n + Rdw <_M Rup_r + Rdw + Rup_r^n$. We reason by contradiction as before. Suppose that $Rup_r^n + Rdw =_M Rup_r + Rdw + Rup_r^n$. Then from the fact that $Rup_r + Rdw + Rup_r^n =_M Rup_r + Rdw + Rup^n$, we have that $Rup_r^n + Rdw =_M Rup_r + Rdw + Rup^n$. On the other hand from the fact that $Rup_r^n + Rdw =_M Rup^n + Rdw$ and that $Rup^n + Rdw <_M Rup_r + Rdw + Rup^n$, we have that $Rup_r^n + Rdw <_M Rup_r + Rdw + Rup^n$ and therefore that $Rup_r^n + Rdw \neq_M Rup_r + Rdw + Rup_r^n$. This contradicts what we have proved before. We conclude therefore that $Rup_r^n + Rdw \neq_M Rup_r + Rdw + Rup_r^n$. ■

Further relations between sets of reflection rules can be inferred from the structure of the graph in Figure 2. Indeed we can exploit the following properties.

PROPOSITION 4.5

Let X and Y be two relational symbols describing a relation among two sets of reflection rules. Let x stand for D, M , or O . Let $(RR)_1, (RR)_2$ and $(RR)_3$ be sets of reflection rules. Then the following properties hold:

IMPLICATION If $(RR)_1 X (RR)_2$, and Y is reachable from X in graph of Figure 1, then $(RR)_1 Y (RR)_2$.

TRANSITIVITY-1 If $(RR)_1 \leq_x (RR)_2$ and $(RR)_2 \leq_x (RR)_3$, then $(RR)_1 \leq_x (RR)_3$.

TRANSITIVITY-2 If $(RR)_1 <_x (RR)_2$ and $(RR)_2 \leq_x (RR)_3$, then $(RR)_1 <_x (RR)_3$.

IDENTITY $(RR)_1 =_x (RR)$.

SUBSTITUTIVITY If $(RR)_1 =_x (RR)_2$ and $(RR)_2 Y_x (RR)_3$, then $(RR)_1 Y_x (RR)_3$.

These properties follow immediately from the definitions of $=_x, \leq_x$ and $<_x$, with $x \in \{O, M, D\}$. Table 2 links the holding of $=_D, =_O, =_M, <_D, <_O, <_M$ of any two sets of

TABLE 2. Strict ordering among (RR)s

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
	\emptyset	$Rdwr$	$Rupr$	Rdw	Rup	$Rupr$ $Rdwr$	$Rupr^n$ Rdw	$Rupr$ Rdw	$Rupr$ Rdw Rup^n	$Rupr$ Rdw Rdw^n Rup^n	$Rupr$ Rdw Rdw^n	$Rupr$ Rdw Rdw^n	Rup^n Rdw	Rup Rdw	Rup Rdw Rup^n
1 \emptyset	$=D$	$<O$ $=M$ $<D$	$=O$ $<M$ $<D$	$<O$ $=M$ $<D$	$=O$ $<M$ $<D$	$<O$ $<M$ $<D$	$<O$ $<M$ $<D$	$<O$ $<M$ $<D$	$<O$ $<M$ $<D$	$<O$ $<M$ $<D$	$<O$ $<M$ $<D$	$<O$ $<M$ $<D$	$<O$ $<M$ $<D$	$<O$ $<M$ $<D$	$<O$ $<M$ $<D$
2 $Rdwr$	$=M$	$=D$	$<M$	$=O$ $=M$ $<D$	$<M$	$<O$ $<M$ $<D$	$<O$ $<M$ $<D$	$<O$ $<M$ $<D$	$<O$ $<M$ $<D$	$<O$ $<M$ $<D$	$<O$ $<M$ $<D$	$<O$ $<M$ $<D$	$<O$ $<M$ $<D$	$<O$ $<M$ $<D$	$<O$ $<M$ $<D$
3 $Rupr$	$=O$	$<O$	$=D$	$<O$	$=O$ $=M$ $<D$	$<O$ $<M$ $<D$	$<O$ $<M$ $<D$	$<O$ $<M$ $<D$	$<O$ $<M$ $<D$	$<O$ $<M$ $<D$	$<O$ $<M$ $<D$	$<O$ $<M$ $<D$	$<O$ $<M$ $<D$	$<O$ $<M$ $<D$	$<O$ $<M$ $<D$
4 Rdw	$=M$	$=O$ $=M$	$<M$	$=D$	$<M$	$<O$ $<M$ $<D$	$<O$ $<M$ $<D$	$<O$ $<M$ $<D$	$<O$ $<M$ $<D$	$<O$ $<M$ $<D$	$<O$ $<M$ $<D$	$<O$ $<M$ $<D$	$<O$ $<M$ $<D$	$<O$ $<M$ $<D$	$<O$ $<M$ $<D$
5 Rup	$=O$	$<O$	$=O$ $=M$	$<O$	$=D$	$<O$ $<M$	$<O$	$<O$ $<M$	$<O$ $<M$	$<O$ $<M$	$<O$ $<M$	$<O$ $<M$	$<O$	$<O$ $<M$	$<O$ $<M$ $<D$
6 $Rupr$ $Rdwr$						$=D$		$<O$ $<M$ $<D$	$<O$ $<M$ $<D$	$<O$ $<M$ $<D$	$<O$ $<M$ $<D$	$<O$ $<M$ $<D$		$<O$ $<M$ $<D$	$<O$ $<M$ $<D$
7 $Rupr^n$ Rdw							$=D$	$<O$ $<M$ $<D$	$<O$ $<M$ $<D$	$<O$ $<M$ $<D$	$<O$ $<M$ $<D$	$<O$ $<M$ $<D$	$<O$ $=M$ $<D$	$<O$	$<O$ $<M$ $<D$
8 $Rupr$ Rdw								$=D$	$=O$ $<M$ $<D$	$<O$ $<M$ $<D$	$=O$ $<M$ $<D$	$<O$ $<M$ $<D$		$<O$ $<M$ $<D$	$<O$ $<M$ $<D$
9 $Rupr$ Rdw Rup^n									$=O$ $=D$	$<O$	$=O$ $=M$ $<D$	$<O$		$<O$	$<O$ $<M$ $<D$
10 $Rupr$ Rdw Rdw^n										$=D$		$<O$ $<M$ $<D$		$<O$ $<M$ $<D$	$<O$ $<M$ $<D$
11 $Rupr$ Rdw Rup^n								$=O$	$=O$ $=M$		$=D$	$<O$		$<O$	$<O$ $<M$ $<D$
12 $Rupr$ Rdw Rdw^n												$=D$		$=O$ $=M$ $<D$	$=O$ $<M$ $<D$
13 Rup^n Rdw							$=M$	$<O$	$<O$ $<M$	$<O$	$<O$ $<M$ $<D$	$<O$	$=D$	$<O$	$<O$ $<M$ $<D$
14 Rup Rdw												$=O$ $=M$		$=D$	$=O$ $<M$ $<D$
15 Rup Rdw Rup^n												$=O$		$=O$	$=D$

reflection rules. A relational symbol X occurs in row $(RR)_1$ column $(RR)_2$ if and only if $(RR)_1 X (RR)_2$. The results in Table 2 can be derived from the facts explicitly stated in the graph in Figure 2 plus the properties IMPLICATION, TRANSITIVITY-1, TRANSITIVITY-2, IDENTITY and SUBSTITUTIVITY.

EXAMPLE 4.6

In order to show that $Rdw <_M Rupr + Rdwr$, which is not implicitly stated in the graph of Figure 2 by a labelled arc, we reason as follows:

1. $Rdwr <_M Rupr + Rdwr$ is stated in the graph of Figure 2;
2. $Rdw =_M Rdwr$ is stated in the graph of Figure 2;
3. $Rdw <_M Rupr + Rdwr$ is derivable from 1. and 2. by SUBSTITUTIVITY.

THEOREM 4.7 (Completeness of the graph in Figure 2)

The graph in Figure 2 is *complete*. That is, for each two sets of reflection rules $(RR)_1$ and $(RR)_2$ and for each relational symbol $X \in \{=D, =O, =M, <D, <O, <M\}$, $(RR)_1 X (RR)_2$ iff $(RR)_1 X (RR)_2$ is derivable from the facts explicitly stated in the graph in Figure 2, and from IMPLICATION, TRANSITIVITY-1, TRANSITIVITY-2, IDENTITY and SUBSTITUTIVITY.

The proof of Theorem 4.7 can be done by showing that, if $<_X [=X]$ does not occur in the

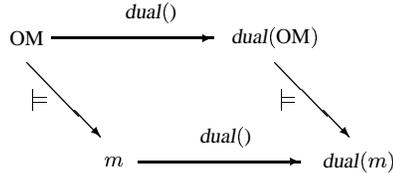


FIGURE 3. A duality graph for satisfiability

slot at row $(RR)_1$ column $(RR)_2$ of Table 2, then there are two OM pairs OM_1 and OM_2 , such that $OM_1 \not\prec_X OM_2$ [$OM_1 \neq_X OM_2$]. The proof, not reported because it is routine and tedious, exploits the same techniques exploited in the proof of Theorem 4.3.

Let us discuss, top to bottom, left to right, some of the results in Table 2. The results in row 1 ($(RR)_1 = \emptyset$) are obvious. It is sufficient to notice that a single reflection up rule does not modify the object theory. Trivially, the dual result holds in the case of reflection down only. The results in rows 2 and 3 have similar explanations. The relations generated by a reflection up rule cannot be compared with those generated by a reflection down rule. A restricted rule generates the same meta and object theories as its corresponding unrestricted rule. It only allows more derivations, all those which (start in one theory and finish in another and) contain at least an application of a reflection rule in presence of open assumptions. The effects of dropping the restriction appear when we have at least one reflection up rule and one reflection down rule. In row 4 the only non-trivial result is in column 6. As from row 3 column 4, Rdw and Rdw_r generate the same meta and object theories. Adding Rup_r to Rdw_r increases the theorems of the metatheory, this proves $Rdw <_M Rup_r + Rdw_r$. To prove that $Rdw <_O Rup_r + Rdw_r$ it is sufficient to consider the case where $\Omega_O = \emptyset$ and $\Omega_M = \bullet("A \vee \neg A") \supset \bullet("B")$. Finally, to prove $Rdw \not\prec_D Rup_r + Rdw_r$ it is sufficient to further consider the fact that any derivation which contains an application of reflection down with open assumptions is not a derivation of the OM pair with $Rup_r + Rdw_r$. In row 5 the result in column 6 can be proved with arguments similar to those for row 4, column 6.

5 A duality property for satisfiability

As already stated in Part I [2], duality principles usually state the preservation of certain logical properties, e.g. provability or satisfiability, under appropriate syntactic transformations on formulas. Examples of duality principles for propositional logic and modal logic can be found in [5] and [1] respectively. In this section we consider satisfiability, in particular we prove (under certain general conditions) the correctness of the graph in Figure 3, where $dual(\cdot)$ is a transformation (defined below) on formulas/ models and OM pairs.

We recall from Part I [2] the notions of dual reflection rule, dual formula, and dual OM pair. See Part I [2] for an explanation of the underlying intuitions.

DEFINITION 5.1 (Dual reflection rule)

For any positive [negative] reflection rule ρ , $dual(\rho)$ is the corresponding negative [positive]

reflection rule.

DEFINITION 5.2 (Dual formula)

For any formula $A \in L_M$, $dual(A)$ is defined accordingly to 1–4:

1. $dual(A) = A$, A does not contain any occurrence of \bullet ;
2. $dual(\bullet("A")) = \neg \bullet(\neg A)$;
3. $dual(A \circ B) = dual(A) \circ dual(B)$, where \circ is either \wedge or \vee or \supset ;
4. $dual(\neg A) = \neg dual(A)$.

Notice that, $dual(\cdot)$ only modifies formulas containing \bullet predicate, and preserves the structure of the connectives. This is not a surprise, as we are interested in a duality principle for metaproperties. To define $dual(\bullet("A"))$ we reason as follows: since $\bullet("A")$ is derivable by Rup from A , $dual(\bullet("A"))$ must be derivable by $dual(Rup)$ (i.e. Rup^n) from A . This implies that $dual(\bullet("A"))$ must be an element of the set of formulas derivable from A by Rup^n . In symbols $dual(\bullet("A")) \in \{\neg \bullet("B") : A \vdash_{\circ} \neg B\}$. We have chosen $dual(\bullet("A")) = \neg \bullet(\neg A)$.

An important property of dualization is that the dual of the dual of an entity (formula, OM pair) is the entity itself (or equivalent, giving some appropriate notion of equivalence). To satisfy this requirement, if A is of the form $\bullet("B")$, $\bullet("B")$ must be equivalent to $\neg \neg \bullet(\neg \neg B)$, which is $dual(dual(A))$. This is not true in all OM pairs. We therefore restrict ourselves only to those OM pairs which satisfy the following property:

$$\vdash_{OM} \bullet("A") \equiv \bullet(\neg \neg A). \quad (5.1)$$

We call the OM pairs which satisfy Condition (5.1), *classical*.

DEFINITION 5.3 (Dual OM pair)

If $OM = \langle O, M, (RR) \rangle$ is classical OM pair, then the *dual* of OM, is the OM pair $dual(OM) = \langle O, dual(M), dual((RR)) \rangle$, where, $dual(M)$, is the theory obtained by substituting the axioms Ω_M of M with $dual(\Omega_M)$ and by adding $\bullet("A") \equiv \bullet(\neg \neg A)$.

Let us define the dualization for the interpretations of the language of the metatheory. We require that an interpretation m satisfies $\bullet("A")$ if and only if the dual interpretation of m satisfies $\neg \bullet(\neg A)$. This implies that:

$$A \in m \iff \neg A \notin dual(m). \quad (5.2)$$

By replacing m with $dual(m)$ and A with $\neg A$ in (5.2) we obtain:

$$\neg A \notin dual(m) \iff \neg \neg A \in dual(dual(m)). \quad (5.3)$$

Combining (5.2) and (5.3) we have that

$$A \in m \iff \neg \neg A \in dual(dual(m)). \quad (5.4)$$

Finally we require that $dual(dual(m))$ satisfies the same formulas as m (this is in order to satisfy the principle that the dual of the dual of an entity is the entity itself). Hence from (5.4) we can conclude:

$$A \in m \iff \neg \neg A \in m. \quad (5.5)$$

Notice that restriction (5.5) is the semantical counterpart of restriction (5.1) for classical OM pairs. In the following, $dual(\cdot)$ is defined only for those interpretations which satisfy condition (5.5). We are now ready to give the notion of dual interpretation.

DEFINITION 5.4 (Dual interpretation)

For each interpretation m of the metalanguage that satisfies condition (5.5), the interpretation $dual(m) = \{A : \neg A \notin m\}$.

A dual interpretation preserves satisfiability.

THEOREM 5.5 (Duality preserves satisfiability)

For each interpretation m of the metalanguage that satisfies condition (5.5), $m \models A$, if and only if $dual(m) \models dual(A)$.

PROOF. We proceed by induction on the complexity of the meta wffs. Let us start by proving that $m \models \bullet("A")$, if and only if $dual(m) \models \neg \bullet("\neg A")$. If $m \models \bullet("A")$, then $A \in m$. By condition (5.5) $\neg \neg A \in m$. By definition of dual interpretation $\neg A \notin dual(m)$ and therefore $dual(m) \models \neg \bullet("\neg A")$. Vice versa, if $dual(m) \models \neg \bullet("\neg A")$, then $\neg A \notin dual(m)$, then, by definition of dual interpretation, $\neg \neg A \in m$. Condition (5.5) implies that $A \in m$ and therefore $m \models \bullet("A")$. The step cases are routine as $dual(\cdot)$ distributes over the connectives. ■

Part I [2] shows that, for classical OM pairs, duality also preserves derivability, that is that:

$$\Gamma_O, \Gamma_M \vdash_{OM} A_O \iff \Gamma_O, dual(\Gamma_M) \vdash_{dual(OM)} A_O \quad (5.6)$$

$$\Gamma_O, \Gamma_M \vdash_{OM} A_M \iff \Gamma_O, dual(\Gamma_M) \vdash_{dual(OM)} dual(A_M) \quad (5.7)$$

where Γ_O, A_O and Γ_M, A_M are a set of object wffs and a set of meta wffs respectively. Combining this latter result with Theorem 5.5 we obtain that the class of models of the metatheory generated by the OM pair $dual(OM)$ is the dual of the class of models of OM. More formally:

THEOREM 5.6

For any classical OM pair OM, an interpretation m is a model of the metatheory of OM if and only if $dual(m)$ is a model of the metatheory of $dual(OM)$.

PROOF. Let m be a model of OM and A a theorem of the metatheory of $dual(OM)$. By [2, Theorem 5.5], $dual(A)$ is a theorem of OM and therefore $m \models dual(A)$. By (5.7), $dual(m) \models dual(dual(A))$. The fact that $dual(dual(A))$ is equivalent to A ensures that $dual(m) \models A$. ■

Theorem 5.6 allows us to extend the results about the models of the metatheory proved in the previous section. The generalization can be done for the set of reflection rules which are the dual of the combinations that generate classical OM pairs. We summarize these results in the following table, constructed from Table 1.

(RR)	Models of $TH_{OM}(M)$
$Rup_r^n + Rdw^n$	The set of formulas consistent with an extension of O
$Rup_r^n + Rdw^n + Rup_r$	The set of formulas consistent with a consistent extension of O
$Rup_r^n + Rdw^n + Rdw_r$	Fixpoint equation dual of (3.10)
$Rup_r^n + Rdw^n + Rup$	The set of formulas consistent with a consistent extension of O
$Rup_r^n + Rdw^n + Rdw$	Any consistent maximal extension of O or the empty set
$Rup^n + Rdw^n$	Any consistent maximal extension of O or the empty set
$Rup^n + Rdw^n + Rup_r$	Any consistent maximal extension of O

6 Conclusion

This paper follows from Part I [2], which investigates the proof theory of OM pairs. The core result of this paper is the definition of a model theory where the models of the metatheory, defined as sets of formulas, are put in correspondence with the object theory, and vice versa. This paper extends the results presented in Part I [2]. In particular, the partial ordering of OM pairs defined in Part I [2] is strengthened to a strict ordering, and the principle of duality for derivability introduced in Part I [2] is extended to satisfiability.

OM pairs are our proposed formalism for studying the foundations of metareasoning. OM pairs separately represent the three main factors of metalevel reasoning (i.e. the object theory, the metatheory, and the reflection principles) using three distinct components, namely, O , M and (RR) . In many approaches a fourth factor is considered, namely, the amalgamation of the object and metatheory. The object theory and the metatheory can be considered as a single theory, or the metatheory can be embedded in the object theory, or the metatheory can contain the object theory. We have shown in Part I [2] that OM pairs can also formalize this aspect by using bridge rules of the form $\frac{M : A}{O : A} M \subseteq O$, and $\frac{O : A}{M : A} O \subseteq M$.

The first advantage of using OM pairs is that it is possible to reconstruct and classify many different and heterogeneous forms of metareasoning in terms of these four parameters. Examples on how different approaches can be reconstructed and compared using OM pairs, are given in the related work section of [2]. To support this comparison, we have studied the proof and model theoretic properties of many possible parameter configurations. For most of the combinations of reflection rules and kinds of object and metatheories (empty, horn, etc.), we have provided an axiomatic characterization, a semantic characterization and a Natural Deduction calculus.

A second advantage is that the explicit representation of these four components allows us to study the effects of each single component on the overall system, for instance, on the theorems or the models of the object theories and metatheories. To this extent, as an example, we have studied the effects of embedding the metatheory into the object theory, when they are connected by different reflection principles. We have proved that amalgamating them does not lead to inconsistent theories, ... and so on.

In order to complete our study of OM pairs, there is at least a further topic that must be accomplished. In these two papers, indeed, we limit ourselves by considering reflection rules on *closed* formulas. There are, however, interesting reflection principles stating the correctness of the metatheory with respect to properties of the object theory on open formulas or terms. Examples of such properties are: *'being an open formula that represents a certain subset of the domain'*, *'being a term of the object theory that denotes a certain individual'*, ... Reflection principles on these properties have been studied in Formal Logics as *non-local reflection principles*, in Meta Logic Programming as *non-ground representations*, while in Knowledge Representation they are often tightly connected to the problem of *quantifying in*. To formalize these more general reflection principles, we have to consider reflection rules on open formulas. OM pairs with reflection rules defined on non-closed formulas will be the subject of further studies. However, it is our opinion that many of the results described in this paper will scale to these new kinds of reflection rules.

Acknowledgements

The research described in this paper owes a lot to the openness and sharing of ideas which exists in the Automated Reasoning Area at ITC-IRST and in the Mechanized Reasoning Group of the University of Trento and Genoa. Discussions with Massimo Benerecetti, Paolo Bouquet, Alessandro Cimatti, Chiara Ghidini, Enrico Giunchiglia, Silvio Imbó and Paolo Traverso, have provided useful insights. We also thank the anonymous reviewer for very useful feedback.

References

- [1] B. F. Chellas. *Modal Logic – an Introduction*. Cambridge University Press, 1980.
- [2] G. Crisculo, F. Giunchiglia, and L. Serafini. A foundation for metareasoning, Part I: The proof theory. *Journal of Logic and Computation*, **12**, 167–208, 2002.
- [3] S. Feferman. Transfinite recursive progressions of axiomatic theories. *Journal of Symbolic Logic*, **27**, 259–316, 1962.
- [4] F. Giunchiglia. Contextual reasoning. *Epistemologia, special issue on I Linguaggi e le Macchine*, XVI:345–364, 1993. Short version in Proceedings IJCAI'93 Workshop on Using Knowledge in its Context, Chambéry, France, pp. 39–49. Also IRST-Technical Report 9211-20, IRST, Trento, Italy, 1993.
- [5] S.C. Kleene. *Introduction to Metamathematics*. North Holland, 1952.
- [6] D. Prawitz. *Natural Deduction - A proof theoretical study*. Almqvist and Wiksell, Stockholm, 1965.

Received 30 October 1997