# The Nature of Reflexive Paradoxes:
# Part I

LEONARD GODDARD and MARK JOHNSTON*

It has been widely recognized that Thomson's "small theorem" [5] is of central importance in understanding the reflexive paradoxes. Several authors (e.g., Herzberger [2], Martin [3], and Goldstein [1]) have exploited it and variations in it in different ways. The purpose of this paper, however, is to show that the formal and philosophical consequences of the theorem are so extensive that they force a general reappraisal of the paradoxes as such. The concern here is to bring out some of these consequences. In Part I there is no intention to promote a particular solution, though in Part II a generalization of Frege's solution is developed. Here, however, the interest is in the general conditions which must be satisfied by any reflexive paradox and any proposed solution. Some of the results which are arrived at are already well known, but they are presented here as interconnected conclusions within a general theory of paradoxicality which arises naturally from Thomson's theorem.

The analysis is carried out entirely in terms of classical two-valued logic since part of the purpose is to discover what can and cannot be done to block the occurrence of reflexive contradictions in a language based on standard quantification theory. We do not want to deny that there are other and perhaps better ways of handling the paradoxes than those which are available in standard two-valued logic, and nothing we say is incompatible with, say,

Martin's use [3] of Thomson's theorem to develop a solution in terms of a logic which admits truth-value gaps; nor, for that matter, is anything we say incompatible in general with the development of other many-valued or para-consistent logics to handle reflexive contradictions. It is our view, however, that since the paradoxes arise in classical two-valued logic, that is where we should begin if we want to understand how and why they arise.

The version of Thomson's theorem which provides a useful starting point is the following thesis of quantification logic:

$\sim\!\alpha$     $\sim(\exists x)(y)(f(y,x) \equiv \sim\!f(y,y))$.

Interestingly, it can be established by a paradox-type of argument from a contradiction which arises in the reflexive case. Thus, as an instance of $(y)A(y) \supset A(x)$ we have $(y)(f(y,x) \equiv \sim\!f(y,y)) \supset (f(x,x) \equiv \sim\!f(x,x))$. The result follows by reductio and generalization on $x$. A similar argument establishes a general version of $\sim\!\alpha$ in the form of a theorem schema: $\sim(\exists x)(y)$ $(A(y,x) \equiv \sim\!A(y,y))$, $x$ free for $y$ in $A$.

Related to $\sim\!\alpha$ we have the following:

$\sim\!\alpha_1$     $\sim(\exists x)(y)(f(y,x) \equiv \sim\!f(x,x))$,

which can be proved along the same lines. This thesis is relevant to some paradoxes, but for the most part the discussion is limited to $\sim\!\alpha$. Most of what is said, however, carries over to $\sim\!\alpha_1$.

Since there is no $x$ which satisfies $(y)(f(y,x) \equiv \sim\!f(y,y))$, there is no unique $x$ which satisfies it. Given the *PM*-theory of descriptions this can be shown formally using $E!(\imath x)A(x) \supset (\exists x)A(x)$. Thus, taking an instance of this along with $\sim\!\alpha$ we have,

$\sim\!\beta$     $\sim\!E!(\imath x)(y)(f(y,x) \equiv \sim\!f(y,y))$.

Similarly,

$\sim\!\beta_1$     $\sim\!E!(\imath x)(y)(f(y,x) \equiv \sim\!f(x,x))$.

In what follows, extensive use is made of the *PM*-theory of descriptions. Except in a few places (notably Part II, Sections 4.4 and 4.5) its special features are not essential but its use in general helps to emphasize many of the points made.

Since $\sim\!\alpha$ is valid, there is no model in which its negation is satisfied, whatever value is given to $f$; i.e., $\alpha$ is false for all substitutions on $f$ in every model. We shall express this by saying that $(\exists x)(y)(f(y,x) \equiv \sim\!f(y,y))$ is a contradiction (or is self-contradictory) for all $f$: i.e., whatever two-place predicate constant is substituted for $f$ in $\alpha$, the result is a self-contradictory sentence. Similar remarks hold for $\beta$ (and for $\alpha_1$ and $\beta_1$). That is, the substitution of a predicate constant for $f$ in $\beta$ results in a self-contradictory sentence.[1]

*1  The nature of simple paradox arguments*     We begin by limiting attention to those paradoxes which seem to depend on direct reflexiveness. These will be called the simple paradoxes. More complicated versions which are cyclic or chain-reflexive will be considered in Section 2.5.

*1.1*   Simple reflexive paradox arguments always involve a contradictory *assumption* of the form $\alpha$ or $\beta$ (depending on the way the argument is presented).

Thus, the paradoxes are typically presented by starting with a postulational definition of an object, class, property, expression, statement, etc. A familiiar example begins: "Consider a catalogue $C$ which lists all and only those catalogues which do not list themselves". Not only is this a contextual definition of $C$, it is also an existential presupposition; and it is from such postulational definitions that the paradox arguments begin. But each of them embodies an instance of $\alpha$ or $\beta$, i.e., a sentence of the form $\alpha$ (or $\beta$) with a predicate constant for $f$. Hence we are never provided with a consistent definition of the postulated item and there is, therefore, a contradiction among the premises of each paradox argument. In fact, we can without loss express every such postulational definition in the form $\alpha$, or alternatively in the form $\beta$. Then the paradox arguments are such that each contains a contradictory premiss of the same form. This is illustrated by the following familiar cases:

(a) *The catalogue.* We suppose that there is a catalogue, say $C$, such that, for any catalogue $y$, $y$ is catalogued in $C$ ($c(y,C)$) iff $y$ is not catalogued in $y$ ($\sim c(y,y)$), i.e., we suppose $(\exists x)(y)(c(y,x) \equiv \sim c(y,y))$. Alternatively, since it is clear from the way the scene is set that uniqueness is presupposed, the supposition is $E!(\imath x)(y)(c(y,x) \equiv \sim c(y,y))$, i.e., $E!C$ where $C$ is $(\imath x)(y)$ $(c(y,x) \equiv \sim c(y,y))$.

(b) *The barber.* We suppose that there is a person, say $B$, such that for any person $y$, $y$ is shaved by $B$ ($s(y,B)$) iff $y$ is not shaved by $y$ ($\sim s(y,y)$); i.e., we suppose $(\exists x)(y)(s(y,x) \equiv \sim s(y,y))$. Alternatively, the supposition is $E!(\imath x)(y)(s(y,x) \equiv \sim s(y,y))$.

(c) *Russell's paradox.* We suppose that there is a class, say $R$, such that for any class $y$, $y$ is a member of $R$ ($\epsilon(y,R)$) iff $y$ is not a member of itself ($\sim\epsilon(y,y)$); i.e., we suppose $(\exists x)(y)(\epsilon(y,x) \equiv \sim\epsilon(y,y))$. Here $\epsilon$ is taken to be a primitive two-place predicate constant not yet characterized by any set-theoretic conditions. Thus the paradox is not being presented as part of set-theory but as part of an applied quantification theory. In particular, there is no assumption that set abstracts are available via an abstraction thesis. On the other hand, given an unrestricted theory of descriptions which permits the formation of a definite description from any predicate constant, we are just as entitled here, as in other cases, to represent the supposition in the stronger form, $E!(\imath x)(y)(\epsilon(y,x) \equiv \sim\epsilon(y,y))$, or to define $R$ by $(\imath x)(y)(\epsilon(y,x) \equiv \sim\epsilon(y,y))$ and affirm $E!R$.

(d) *Grelling's paradox.* We suppose that there is an adjective, *Het*, such that for any adjective $y$, $y$ has the feature described by *Het* ($d(y,Het)$) iff $y$ does not have the feature described by $y$ ($\sim d(y,y)$); i.e., we suppose $(\exists x)(y)$ $(d(y,x) \equiv \sim d(y,y))$, alternatively $E!(\imath x)(y)(d(y,x) \equiv \sim d(y,y))$. Here quotation marks are omitted for simplicity. They can be introduced by regarding $d(y,x)$ as an abbreviation for $d'(qu(y), qu(x))$, where $qu$ is a quotation function. Alternatively, taking the *is* of predication to be a two-place relation, a different formulation of the presupposition is $(\exists x)(y)(y \; is' \; x \equiv \sim(y \; is' \; y))$

where $y$ *is'* $x$ is an abbreviation for $qu(y)$ *is* $x$. These last two remarks show what is of course well known: that the mere introduction of linguistic levels via quotation marks is not enough to remove Grelling's paradox.

(e) *The liar*. It may seem that the simple liar, 'This statement is false', fails to fit the pattern of the other paradoxes both because 'false' is a one-place predicate and because no quantification is involved. However, we can obtain an equivalent of the simple liar as follows. We suppose that one and only one statement $S$ is made in a particular spatio-temporal stretch $t$, namely the statement that all statements made in $t$ are false. We suppose, too, that $S$ is either true or false. These two suppositions are equivalent to the assumption that for any statement $y$ made in $t$, $S$ is a true statement about $y$ $(tr(y,S))$ iff $y$ is a false (not true) statement about $y$ $(\sim tr(y,y))$, since $S$ is the only statement made in $t$. Hence, where the variables range over statements made in $t$, the assumption is $(\exists x)(y)(tr(y,x) \equiv \sim tr(y,y))$; alternatively, $E!(\imath x)(y)(tr(y,x) \equiv \sim tr(y,y))$. Here, since we are taking it that only one statement is made in $t$, the quantifiers are doing no real work and the assumption is effectively the degenerate case: $tr(S,S) \equiv \sim tr(S,S)$, i.e., $S$ is a true statement about $S$ iff $S$ is a false statement about $S$. Thus we have an analogue of the simple liar in the form: $S$ is the statement that $S$ is false, since this too is equivalent to the assumption that $S$ is a true statement about $S$ iff $S$ is a false statement about $S$.

Since the quantifiers are idle (not vacuous) in $(\exists x)(y)(tr(y,x) \equiv \sim tr(y,y))$ we could equally well have developed the simple liar as a case of $\alpha_1$, i.e., $(\exists x)(y)(tr(y,x) \equiv \sim tr(x,x))$. And in fact it is interesting to see that this version of $\alpha_1$, when interpreted over a suitable domain in which the quantifiers are not idle, generates a variant of the Cretan paradox: e.g., let $S$ be as before but suppose that other statements are made in $t$ all of which are false. We call this the strict Cretan since it differs from the ordinary Cretan in the assumption that all other statements referred to by $S$ are in fact false.

A mixed Grelling-liar paradox can be developed from a nondegenerate case of $\alpha$: $(\exists x)(y)(tr(y,x) \equiv \sim tr(y,y))$. Thus, consider a set of type sentences with members of the sort 'This is a . . . sentence' where the blanks are filled by suitable adjectives or adjectival phrases, e.g., 'This is a long sentence', 'This is a printed-in-red-ink sentence'. Allow that any member of the set can be used to refer to any member of the set (including itself). On each occasion of its use a given sentence will yield a true statement or a false statement (assuming bivalence) depending on the features of the sentence referred to. Thus, given some criterion for being a long sentence (contains more than five words, say), 'This is a long sentence' yields a true statement about 'This is a printed-in-red-ink sentence' and a false statement about 'This is a short sentence'. In particular, then, they yield true or false statements when used to refer to themselves. Call those which yield false statements when used to refer to themselves, self-falsifying sentences. Then consider 'This is a self-falsifying sentence'. This sentence, abbreviated $F$, yields a true statement about any sentence in the set iff the sentence to which it refers yields a false (not true) statement about itself, i.e., for any sentence $y$ in the set, $F$ yields a true statement about $y$ $(tr(y,F))$ iff $y$ does not yield a true statement about $y$ $(\sim tr(y,y))$. Hence, $(\exists x)(y)(tr(y,x) \equiv \sim tr(y,y))$. Here the predicate constant $tr$ relates sentences

rather than statements; as in the Grelling, therefore, we take it that the predicate constant absorbs quotation.

The simple liar can thus be seen either as a degenerate case of the strict Cretan, in which $S$ is the only statement made in $t$, or as a degenerate case of the Grelling-liar, in which $F$ is the only element of the domain.

*1.2*   It might be objected at this point that the formalization of the postulational definitions is misleading since in the informal exposition of the examples we have used restricted quantifiers (for all catalogues $y$, for all classes $y$, etc.), but in the formalism we have used unrestricted quantifiers. Hence, it might be said, we have misrepresented the actual form of the assumptions being made. Moreover, it might seem, if the assumptions are represented correctly, it becomes immediately obvious that there is no genuine paradox because there is no contradiction.

Thus, suppose we represent 'there is a catalogue, say $C$, such that, for all catalogues $y$, $c(y,C)$ iff $\sim c(y,y)$' by $(y)(c^*y \ \& \ c^*C. \supset (c(y,C) \equiv \sim c(y,y)))$, where $c^*$ is the predicate ' . . . is a catalogue'. The conclusion from the reflexive case would not then be a contradiction but instead $\sim c^*C$, i.e., $C$ is not a catalogue.

Such a move, however, simply pushes the problem one step back. For there seems to be no reason why we should not limit the universe of discourse to catalogues, and if we make this explicit by affirming $(y)c^*y$, either we recover the contradiction or we have to conclude that there cannot be a domain consisting only of catalogues, and this might be thought to be paradox enough. This point is taken up in Part II, Section 4.6.

In effect, our presentation of the paradoxical assumptions focuses attention on the paradoxes when construed in this limited way over restricted domains. For we have achieved the restriction on the quantifiers by using unrestricted quantifiers over restricted domains (where of course the domain is different in each example) instead of using restricted quantifiers over an unrestricted domain. So construed, the paradoxical assumptions are contradictory.

*1.3*   Paradox arguments are often represented as being argument-pairs of the form:

(a)  $p_0 \vdash \sim p_0$
(b)  $\sim p_0 \vdash p_0$

where $p_0$ is some sentence-constant, e.g., $c(C,C)$. And since we have in general, $p \vdash p$ and $\sim p \vdash \sim p$, we may represent them as:

(a')  $p_0 \vdash p_0 \ \& \ \sim p_0$
(b')  $\sim p_0 \vdash p_0 \ \& \ \sim p_0$.

That is, the assumption of $p_0$ leads to a contradiction and so does the assumption of $\sim p_0$. It should now be clear, however, that this is a misrepresentation. The assumption $c(C,C)$, say, does not yield $\sim c(C,C)$ unless it is combined with an appropriate instantiation case of $(\exists x)(y)(c(y,x) \equiv \sim c(y,y))$, viz., $c(C,C) \equiv \sim c(C,C)$, or of $E!(\imath x)(y)(c(y,x) \equiv \sim c(y,y))$. And it is just this further premiss

which is presupposed when it is said (as part of the lead-in to the argument) that $C$ is the catalogue of all and only those catalogues which do not list themselves. Hence $(\exists x)(y)(c(y,x) \equiv \sim c(y,y))$, or $E!(\imath x)(y)(c(y,x) \equiv \sim c(y,y))$, is an essential premiss in *each half* of the paradox argument.

In general, since an $\alpha$-instance (or a $\beta$-instance) is *assumed* before a paradox argument can begin, it is always an essential premiss in both parts (a) and (b). Thus, in terms of an $\alpha$-instance, the real structure of the arguments for a given predicate constant $f_0$ and individual constant $x_0$ is:

(a'')   $(\exists x)(y)(f_0(y,x) \equiv \sim f_0(y,y))$                     Ass(1)
        $(y)(f_0(y,x_0) \equiv \sim f_0(y,y))$                                EI
        $f_0(x_0,x_0) \equiv \sim f_0(x_0,x_0)$                               UI
        $f_0(x_0,x_0)$                                                         Ass(2)
        $\sim f_0(x_0,x_0)$

(b'')   $(\exists x)(y)(f_0(y,x) \equiv \sim f_0(y,y))$                     Ass(1)
        $(y)(f_0(y,x_0) \equiv \sim f_0(y,y))$                                EI
        $f_0(x_0,x_0) \equiv \sim f_0(x_0,x_0)$                               UI
        $\sim f_0(x_0,x_0)$                                                   Ass(2)
        $f_0(x_0,x_0)$

Alternatively, if we start with the $\beta$-instance $E!(\imath x)(y)(f_0(y,x) \equiv \sim f_0(y,y))$ then it becomes Ass(1), the $\alpha$-instance is an immediate consequence, and the argument continues as above. In fact, of course, except as part of a parlour-game demonstration, the last two steps (which are usually taken to be the most significant) are unnecessary, i.e., we do not need to affirm Ass(2). Then, the two halves of the paradox argument are identical.

*1.4*   We have presented the argument against a background in which it is already known that $\alpha$ and $\beta$ are contradictions, and in doing so we have of course lost the paradox since it is in no sense paradoxical to argue validly from a contradiction to a contradiction. But even in a context in which this is not known, it is not obvious where the paradox lies. If the argument is set out formally with the $\alpha$- or $\beta$-instance made explicit (and there is no argument at all if it is not), then we have a simple reductio which establishes the $\alpha$- or $\beta$-instance as contradictory; and there is nothing paradoxical about that. More generally, if the argument is set out with parameter $f$ standing in for any two-place predicate constant, then we simply have a formal proof of $\sim\alpha$ ($\sim\beta$) which differs from the proof given earlier only in its manner of presentation. In general, if what is thought to be characteristic of a paradox is that there be some sentence $p$ which *of itself*, without being conjoined to further assumptions (other than logical theses), yields $\sim p$, and $\sim p$ *of itself* yields $p$, then the simple paradoxes fail to supply us with a paradox. However, we do not wish to quibble about nomenclature. What remains puzzling perhaps, if not paradoxical, is that $\alpha$ and $\beta$ *are* contradictions. Why is it that there *cannot be* a barber, a catalogue, a class etc., of the kind defined? We take up this point in Section 2 below.

*2 The nature of the contradiction*     To have seen what is wrong with the paradox arguments is not yet to have explained why what is going wrong is

going wrong. For what is puzzling about the arguments is not that we arrive at a contradictory conclusion but that we should ever find ourselves in the position of wanting to affirm a contradictory premiss. The situation is quite different, for example, from that which would arise if we were to affirm that there is a person, say Jack, who loves everyone, and at the same time another person, say Jill, whom no one loves. Here the contradiction is obvious and no concern would arise from the fact that we could validly derive the conclusion that Jack loves Jill and Jack does not love Jill. Yet formally speaking there is no difference between this argument and the simple paradox arguments. Each begins with a negated instance of a quantification thesis and then proceeds by standard moves to a more explicit contradiction expressed in terms of the postulated individuals. In this particular case the initial assumption is the negation of the thesis $(\exists x)(y)f(x,y) \supset (y)(\exists x)f(x,y)$ for a given predicate constant $f_0$ (loves): i.e., pushing the negation through, the assumption is $(\exists x)(y)f_0(x,y)$ & $(\exists y)(x)\sim f_0(x,y)$. So (skipping a few steps to get the parallel with the paradox arguments) we then have, essentially by EI, $(y)f_0(x_0,y)$ & $(x)\sim f_0(x,y_0)$; and by UI (etc.), $f_0(x_0,y_0)$ & $\sim f_0(x_0,y_0)$.

The difference between this argument and the simple paradox argument lies solely in our attitude toward the initial assumptions. There is no temptation to think that there are a Jack and Jill who satisfy the postulated condition. By contrast, there is a temptation to think that there are such individuals as $C$, $B$, $R$, etc., which satisfy the postulated $\alpha$-instances. That is, there is a temptation to think that the $\alpha$-instances *ought* to be true in spite of the fact that they are contradictory. This is what makes for the feeling of paradox. And the main reason for this temptation is the belief that the postulated $\alpha$-condition $(y)(f(y,I) \equiv \sim f(y,y))$, where $I$ is the item chosen to instantiate the existential quantifier, fails only in the one special case where the condition is applied to $I$ itself to yield $f(I,I) \equiv \sim f(I,I)$; for all other instantiations it is true. Thus, it seems, given the exclusion of the reflexive case, the $\alpha$-condition would be acceptable.

But this belief is false. For although it is true that the formal inconsistency arises from the reflexive case (this is discussed in more detail in Section 2.3), it is not true in every case that the $\alpha$-instance is acceptable even if some way could be found of preventing its application to the particular item $I$. In fact the familiar paradoxes, which carry an initial air of plausibility, form a very special class of cases. This becomes important when one considers what would count as an acceptable solution. To show this, we first look at some other examples.

*2.1*   Since $\alpha$ is a contradiction for all $f$ in every domain, it follows that we can construct postulational definitions like those in Section 1.1 using *any* two-place predicate constant. Indeed, we do not need to restrict ourselves to two-place predicates. Any instance of the generalized form of $\alpha$, i.e., $(\exists x)(y)(A(y,x) \equiv \sim A(y,y))$, will yield a paradox. However, not all members of the family have the same air of paradoxicality as the familiar ones, though some do.

Thus, if we take the predicate 'is created by' ($cr$) and suppose that there is an individual, say $G$, such that, for any individual $y$, $y$ is created by $G$ ($cr(y,G)$) iff $y$ is not created by $y$ ($\sim cr(y,y)$), i.e., $(\exists x)(y)(cr(y,x) \equiv \sim cr(y,y))$, then we

have something like a standard paradox. Putting it in a more familiar form we may say: Let $G$ be the creator of all and only those who do not create themselves. Who creates $G$?

By contrast, take the predicate constant 'to the right of' ($r$) and postulate that there is an individual, say $L$, such that for any individual $y$, $y$ is to the right of $L$ ($r(y,L)$)) iff $y$ is not to the right of itself ($\sim r(y,y)$)), i.e., $(\exists x)(y)$ $(r(y,x) \equiv \sim r(y,y))$. Here we may seem to have a paradox. We may say: Let $L$ be an object such that anything is to the right of $L$ iff it is not to the right of itself. Is $L$ to the right of itself or not? But the paradox lacks plausibility for the very obvious reason that since it is true of every individual that it is not to the right of itself, the postulational definition effectively asserts that there is some individual such that everything is to the right of it. That is, $L$ is postulated to be an ultimate left-hand object. But our use of 'to the right of' is incompatible with there being such an object.[2] Thus, the postulational definition can be dismissed as false by meaning even if it is not recognized to be formally inconsistent, though of course it is formally inconsistent quite independently of the meaning of the predicate constant and for exactly the same reason as all cases of $\alpha$ are. Similarly, if we take the predicate constant 'less than' and restrict the domain to natural numbers, then, since every number is not less than itself, the appropriate $\alpha$-instance $(\exists x)(y)((y < x) \equiv \sim(y < y))$ in effect affirms the existence of a greatest number $N$. Hence it can be dismissed as false by meaning independently of the paradoxical case which arises when we ask whether or not $N$ is less than itself.

There is, then, a distinction to be made between those instances of $\alpha$ which, were it not for the formal inconsistency arising in one particular case, would be acceptable, and those which remain unacceptable even if the formal inconsistency could be removed. The former we call plausible (for want of a better word), the latter implausible. We can make this distinction more precise as follows: For a given $\alpha$-instance $\alpha_0$ expressed in terms of a predicate constant $f_0$, let $M(f_0)$ be a set of meaning postulates on $f_0$. Then, $\alpha_0$ is implausible if $\{M(f_0), \alpha_0\}$ is inconsistent independently of the reflexive case. If this condition is satisfied, we shall say that $\alpha_0$ is false by meaning as well as being self-contradictory by virtue of the reflexive case. That is, the meaning of the predicate constant precludes the existence of the postulated individual independently of the nonexistence entailed by the contradiction which arises in the reflexive case.

Of course, the applicability of this condition will depend on there being suitable meaning postulates. In some cases, however, there is no difficulty in applying it. For example, let $\alpha_0$ be $(\exists x)(y)((y < x) \equiv \sim(y < y))$, where the domain is the natural numbers, and let $M(<)$ be (1) $(x)(y)(y < x \supset \sim(x < y))$, (2) $(x)(x < x + 1)$. Then a consequence of $\alpha_0$ and $(y) \sim (y < y)$ (derivable from (1)), is $N + 1 < N$, where $N$ is an arbitrary constant chosen to instantiate the existential quantifier; but a consequence of (1) and (2) is $\sim(N + 1 < N)$. Hence $\{M(<), \alpha_0\}$ is inconsistent independently of the reflexive case which arises directly from $\alpha_0$ itself.

Here, as in other cases, the meaning postulates characterize formal features of the relation designated by the predicate constant. It should not be thought, however, that there is some common formal condition which

characterizes the relation in every plausible (or implausible) paradox and which could therefore be used as a uniform criterion for distinguishing the plausible paradoxes from others. Thus, in the standard examples, the relations in question are usually taken to be nonreflexive (some classes are members of themselves, some not; some catalogues list themselves, some do not; and so on). And even in the creator paradox, there is an implicit presupposition, which seems to carry plausibility, that although most beings do not create themselves, there is at least one self-creating being, namely $G$. So here, too, the relation is nonreflexive. By contrast, in the implausible to-the-right-of and less-than paradoxes, the relations are irreflexive. Moreover, it is possible to construct equally implausible examples using reflexive relations (e.g., equality over the natural numbers). So it might perhaps be thought that all and only nonreflexive relations give rise to plausible paradoxes. However, such a basis for the distinction cannot be maintained. Given the nonreflexive relation 'loves', the paradoxical assumption that all and only those who do not love themselves love the same one individual, $K$, is entirely implausible. At the same time, suppose no catalogues list themselves, so that $c$ is an irreflexive relation; it nevertheless seems initially plausible to suppose that we could compile a catalogue of all catalogues which do not list themselves, i.e., a catalogue of all catalogues. Or again, suppose that no class is a member of itself; it seems plausible to suppose that we could form a class of all such classes, i.e., a class of all classes.

*2.2*  There might, then, be some grounds for distinguishing what might be called semantic paradoxes from others by classifying as semantic those which can be dismissed as false by meaning. But this distinction is quite different from Ramsey's distinction [4] between logical and semantic paradoxes. The basis of Ramsey's classification is between those paradoxes which can be set up entirely within a formal system (in fact set theory) and those which require additional (semantic) concepts. Or we may say, perhaps, that the logical paradoxes for Ramsey are those such that the domain over which the quantifiers range contains only formal objects (numbers, classes, etc.), while the domain in the case of the semantic paradoxes contains nonformal objects such as people, catalogues, statements, etc. But this seems to be the least important feature of the paradoxes, even though Ramsey's classification had the advantage of simplifying the type-theoretical solution of the "logical" paradoxes. What is important about the paradoxes is their common feature, not their differences. This common feature is that the presuppositions in each of them have the same formal structure. And it is the formal structure of the presupposition, not its interpretation over a domain, which lies at the heart of the paradoxes. Each paradox presupposes an α-instance which is demonstrably false in pure quantification theory interpreted over *any* domain. In this sense, all are logical; for the kinds of objects admitted to the domain are irrelevant to the demonstration. This applies to the implausible examples no less than to those which are plausible.

Ramsey's distinction is therefore misleading. And it would be equally misleading to argue in the converse direction from the fact that some paradoxes can be removed by one kind of known solution while others cannot (e.g.,

Russell's can be removed by simple type theory while Grelling's requires order theory), to the conclusion that the paradoxes should be classified differently when they yield to different solutions, for this, too, ignores the similarity of formal structure in the different paradoxes.

A different way in which the paradoxes might be thought to be distinguishable arises from the fact that we may feel intuitively that the existential presupposition is reasonable in some cases but not in others. Thus, it might be thought, we have good reason to suppose that $C$ exists since we can make a start on compiling it, or that $R$ exists since we can make a start on constructing it—we include the class of men but omit the class of abstract items. By contrast, we have no good reason to suppose that the barber exists—the supposition that he does is simply an ad hoc assumption, and what the paradox shows is that the assumption is false. But this way of distinguishing the paradoxes also fails to acknowledge their common formal structure, especially so if there is an intended covert implication that the existential presupposition in the case of *The Barber* is contingently false. Given an appropriate domain (men), a relation defined on that domain (shaves), and a predicate constant which designates that relation ($s$), we can construct the description $(\imath x)(y)$ $(s(y,x) \equiv {\sim}s(y,y))$, and it is always reasonable to ask whether or not there is an individual satisfying a given description whatever initial intuitions one might have about the answer. The fact that we can then go on to prove (by way of a paradox argument) that there is no individual answering to the description, i.e., that the barber does not exist, no doubt confirms our intuitions about the outcome, but that is irrelevant to the logic of the outcome. Logically speaking, *The Barber* is identical with the other paradoxes. *Necessarily*, there is no barber, no catalogue $C$, no Russell class, no adjective with the sense intended for *Het*, no statement yielded by the liar sentence,[3] no ultimate left-hand object, and no creator satisfying the description $G$: because there *cannot be* in *any* domain *any* individual satisfying a description of the form $(\imath x)(y)(f(y,x) \equiv {\sim}f(y,y))$. What we must conclude, then, is that our original intuitions about the reasonableness of positing the existence of $C$ and $R$ were mistaken. It was just as unreasonable (or as reasonable) to make those assumptions as it was to assume that the barber exists. For whatever catalogue we think we are compiling when we begin to list those catalogues which do not list themselves, we are not compiling $C$ and we are not even making a start on compiling $C$, since $C$ cannot exist. The catalogue we begin to compile has some entries in common with those which $C$ was intended to have, but it is still not $C$ (see Part II, Sections 4.5 and 4.6). Similarly, whatever class we begin to construct when we "collect" together those classes which are not members of themselves, it cannot be $R$. This is what the paradox arguments show, or more generally, what $\beta$ shows. No doubt we have received a surprise, but surprises cannot be the basis for making logical distinctions; and there is no logical basis for making distinctions. On the contrary, the similarity of formal structure forces us to treat all paradoxes as equals. Each exhibits a formal contradiction of the same form.

This similarity of structure is fundamental. It justifies the intuitive view that piecemeal solutions are not solutions at all and that there must be a single solution which applies to every paradox in the same way. In view of our

distinction between semantic paradoxes and others, however, this intuition has to be interpreted with caution. What is true is that since every paradoxical assumption has the same formal structure, any technique which is devised to remove the formal inconsistency arising from that structure must apply equally to all cases. But it should not be expected that the mere removal of the inconsistency will always result in a true assumption, for the modified assumption, though no longer formally inconsistent, may remain false by meaning. That the modified assumption should be true can only be a requirement for that limited class of cases which are initially plausible. Nevertheless, a uniform technique for removing the formal inconsistency would be a general resolution of the paradoxicality of the paradoxes. The question is, however, whether there is such a technique; but that question cannot be answered unless we know exactly why the contradiction arises.

*2.3*    Why is $(\exists x)(y)(f(y,x) \equiv \sim f(y,y))$ contradictory?

Since $f(y,x) \equiv \sim f(y,y)$ is equivalent to $\sim(f(y,x) \equiv f(y,y))$, the formula which is the scope of the innermost quantifier is the negation of a condition for the equality of $x$ and $y$. Thus, from the identity schema $(y = x) \supset (A(y,x) \equiv A(y,y))$, we have $(y = x) \supset (f(y,x) \equiv f(y,y))$; hence $\sim(f(y,x) \equiv f(y,y)) \supset (y \neq x)$, for all $x$ and $y$: i.e., $(x)(y) ((f(y,x) \equiv \sim f(y,y)) \supset (y \neq x))$ is a thesis. At the same time, if we take an instantiation case of $\alpha$, say $(y)(f(y,I) \equiv \sim f(y,y))$, there is nothing which prevents us from taking $I$ to be a value for $y$, indeed we are *required* to include $I$ in the values over which the unrestricted universal quantifier ranges. But that is to take $y = x$ for that particular instantiation. Thus, we have a condition which entails $y \neq x$, for all $x$ and $y$, yet the quantification over that condition requires us to include the case for which $y = x$ for some $x$ and $y$.

*2.4*    The contradiction arises, therefore, simply because of the incompatibility between the quantificational structure of $\alpha$ and the implications of its matrix (the scope of the innermost quantifier). The arrangement of quantifiers permits an instantiation case which presupposes an equality the negation of which is entailed by the matrix. The quantificational structure and the matrix are each innocuous in other contexts, but together they constitute a sufficient condition for contradiction in the reflexive case.

We can thus describe a general form of the simple reflexive paradoxes as follows:

C    Given a formula $A(y,x)$ containing just two free variables $x$ and $y$ and no bound variables, the formula $(\exists x)(y)A(y,x)$ will be a formal contradiction, and in particular a reflexive contradiction, in case $A(y,x) \supset (y \neq x)$ is a thesis.

For since in general we have $(y)A(y,x) \supset A(x,x)$, subject to the usual proviso, and since $A(x,x) \supset (x \neq x)$ follows from $A(y,x) \supset (y \neq x)$, we conclude $\sim(y)A(y,x)$, i.e., universally quantifying w.r.t. $x$, $\sim(\exists x)(y)A(y,x)$. Thus $\vdash A(y,x) \supset (y \neq x)$ is a sufficient condition for a formula to be a reflexive contradiction in the special case in which it takes the form $(\exists x)(y)A(y,x)$.

It should be noted, however, that this condition introduces a more general notion of a reflexive contradiction. For since $A(x,y)$ need not be in the form

of an equivalence, we have lost the familiar equivalence structure which characterizes the simple paradoxes. Thus the trivial case of a reflexive contradiction is now given by $(\exists x)(y)(y \neq x)$; at the same time, formulas more complex than those which give rise to the simple paradoxes get classified as reflexive contradictions. Hence the simple paradoxes, whether plausible or implausible, cover only a small band of the spectrum determined by C.

It is clear, too, that the restriction to two variables, though a characteristic of the simple paradoxes, is independent of the condition which gives rise to the reflexive contradiction. This condition, that the quantificational structure of a formula be incompatible (in the way described) with the implications of its matrix, does not of itself impose any limitations on the number of variables or quantifiers. For example, $(\exists x)(y)(z)(y = x \,\&\, z \neq x)$ is a reflexive contradiction since the matrix entails $y \neq z$ for all $y$, $z$, but the quantifier arrangement permits the same value to be chosen to instantiate the universal quantifiers over $y$ and $z$, i.e., permits $y = z$ for arbitrarily chosen $y$ (independently of the value chosen to instantiate the existential quantifier). Hence as a permissible instantiation case we have the contradiction $y_1 = x_1 \,\&\, y_1 \neq x_1$. Or again, given the identity condition $(x = z) \supset (B(y,x) \equiv B(y,z))$, we have $(B(y,x) \equiv {\sim}B(y,z)) \supset (x \neq z)$. Consequently, $(\exists x)(y)(z)(B(y,x) \equiv {\sim}B(y,z))$ is a reflexive contradiction; in particular, then, $(\exists x)(y)(z)(y \,\epsilon\, x \equiv {\sim}(y \,\epsilon\, z))$[4] is "paradoxical". More generally, $(\exists x)(y)(z)A(x,y,z)$ will be a reflexive contradiction in case any of $A(x,y,z) \supset (y \neq x)$, $A(x,y,z) \supset (z \neq x)$, or $A(x,y,z) \supset (y \neq z)$ are theses.

We may therefore formulate a general condition for a reflexive contradiction since an incompatibility of the kind exhibited by the simple case will arise whenever a formula with initial quantifiers is such that the formula which is the scope of the innermost quantifier entails the negation of an equality condition which is presupposed by the quantifier arrangement, in the sense that permissible instantiation cases of the formula commit us to that equality condition in particular cases, i.e.,

$\text{C}_1$     Given a formula of the form $QA$, where $Q$ is a string of quantifiers, a reflexive contradiction will arise if $A$ entails the negation of an equality condition or conditions and permissible instantiation cases of $QA$ presuppose any such condition(s).

Here, $A$ is not precluded from containing other quantifiers or, for that matter, free variables which are not bound by quantifiers in the string $Q$. To impose restrictions on $A$ by limiting the number and kind of variables it may contain would be to limit the generality of the condition in a way which is not required by the intuitive idea of incompatibility which we are trying to express. However, $\text{C}_1$ achieves generality at the cost of vagueness since no precise formulation of the relevant equality conditions is specified and in fact cannot be specified if $A$ is allowed to contain bound variables, or free variables which are not bound in $QA$. In such cases the particular inequalities which are incompatible with permissible instantiation cases of $QA$ may not be simple inequalities between two distinct variables but instead involve a choice operator, and the general inequality condition entailed by $A$ may not be a single inequality but a complex formula involving several inequalities. This is illustrated in the

next section where it is shown that a relatively simple formula $QA$, where $A$ contains just one bound variable, is such that $A$ entails a quite complex inequality condition. In general, the actual form of the inequality conditions depends crucially on the variables (free or bound) which occur in $A$ but not in $Q$ and it is for this reason that the general criterion $C_1$ is inevitably vague. An equivalent criterion $C_3$ which avoids these complexities is given in Part II, Section 1.2.

$C_1$ expresses a sufficient condition for any formula to be a reflexive contradiction: but unlike C, sufficiency cannot be proved formally because of the inherent vagueness of $C_1$. This problem will therefore be left until $C_3$ has been formulated; so too will the question of necessity. We note, however, that although $C_1$ is not expressed in precise formal terms, it is nevertheless a purely formal condition in the sense that it is independent of the meaning of any predicate constants which may occur in $A$ and independent of the choice of domain. Given that a formula $QA$ satisfies $C_1$, then it is a reflexive contradiction however it is interpreted.

*2.5*   It is clear, then, that the contradiction $\alpha$ and instances of it which are associated with the familiar paradoxes are simply special cases of formulas which satisfy $C_1$. What the general condition shows, therefore, is that the simple paradoxes are not the only reflexive contradictions. We have already given some obvious examples to illustrate this, but we now show that the cyclic paradoxes are identified as reflexive contradictions in terms of $C_1$. That is,

$$\alpha_2 \qquad (\exists x)(y)(f(y,x) \equiv \sim(\exists z_1)(\exists z_2) \ldots (\exists z_n)(f(y,z_1) \& f(z_1,z_2)$$
$$\& \ldots \& f(z_n,y))$$

is a reflexive contradiction, as defined above, and consequently $\sim\alpha_2$ is a thesis of quantification theory. As a special case we then have the extended class paradox which arises from the assumption,

$$(\exists x)(y)(y \in x \equiv \sim(\exists z_1) \ldots (\exists z_n)(y \in z_1 \& z_1 \in z_2 \& \ldots \& z_n \in y)).$$

To show which inequalities are involved, we restrict attention to the one-cycle case,

$$\alpha_2^1 \qquad (\exists x)(y)(f(y,x) \equiv \sim(\exists z)(f(y,z) \& f(z,y))).$$

The argument is easily extended to the $n$-cycle case.

The most direct way of discovering the inequalities, so demonstrating that $\alpha_2^1$ is a reflexive contradiction, is to see what instantiations need to be made to carry the paradox argument through.

Taking $x_0$ as an instantiation value for the initial existential quantifier in $\alpha_2^1$, we have:

(i)    $(y)(f(y,x_0) \equiv \sim(\exists z)(f(y,z) \& f(z,y)))$.

So, instantiating the universal quantifier with $x_0$, i.e. *by taking* $y = x_0$,

(ii)    $f(x_0,x_0) \equiv \sim(\exists z)(f(x_0,z) \& f(z,x_0)))$.

But by standard laws,

$$f(x_0,x_0) \supset (\exists z)(f(x_0,z) \& f(z,x_0)).$$

Hence,

(iii)   $\sim f(x_0,x_0)$.

From (ii) and (iii), therefore,

(iv)   $(\exists z)(f(x_0,z) \& f(z,x_0))$.

Now take $z_0$ as an instantiation for the existential quantifier. Then,

(v)   $f(x_0,z_0) \& f(z_0,x_0)$.

Hence,

(vi)   $f(z_0,x_0)$.

But from (i), instantiating the universal quantifier with $z_0$, i.e., *by taking* $y = z_0$, we have,

(vii)   $f(z_0,x_0) \equiv \sim(\exists z)(f(z_0,z) \& f(z,z_0))$.

So, from (vi) and (vii),

(viii)   $\sim(\exists z)(f(z_0,z) \& f(z,z_0))$.

But from (v) we have:

(ix)   $(\exists z)(f(z_0,z) \& f(z,z_0))$.

For the argument to go through, then, we have to take $y = x_0$ and $y = z_0$, where $z_0$ is an arbitrary individual satisfying the condition $f(x_0,z) \& f(z,x_0)$. Hence we should expect the scope of the innermost of the initial quantifiers in $\alpha_2^1$ to entail not both $y = x$ and $y = (\varepsilon z)(f(x,z) \& f(z,x))$, where $\varepsilon$ is the choice operator. And in fact it is a straightforward matter to show:

$$(f(y,x) \equiv \sim(\exists z)(f(y,z) \& f(z,y))) \supset ((y \neq x) \vee y \neq (\varepsilon z)(f(x,z) \& f(z,x))).$$

Here we have arrived at the condition which establishes $\alpha_2^1$ as a reflexive contradiction in terms of $C_1$ by tracking through the standard paradox argument, but in this respect the paradox argument was used simply as a heuristic device. Given that the condition can be demonstrated directly, and independently of the paradox argument, as it can, we thereby establish $\sim\alpha_2^1$ as a thesis of quantification theory. Hence the paradox argument begins from a contradiction and simply proceeds to a more explicit contradiction by utilizing instantiations which can be "read off" the condition. Exactly this is true of every reflexive paradox.

*3 The nature of solutions*     It is not always clear what is being asked for when a solution of the paradoxes is demanded. Often, the expectation has been that whatever technique is employed to remove the formal inconsistency should result in a modified sentence which is true. But this conflates the problem of paradoxicality with that of plausibility. This is not surprising, perhaps, since the familiar paradoxes are all such that the formal inconsistency arising from the reflexive case seems to be (or is generally assumed to be) the

only false case; and if this is so, the removal of the contradiction does result in a true sentence. At the more general level, however, plausibility and para-doxicality have to be kept distinct.

*3.1*   Any technique which successfully removes the contradiction, whether or not it results in a true sentence, we call a resolution rather than a solution. Restricting attention to the two-variable case for the moment, it should now be clear that a resolution can only be achieved if the value chosen to instantiate the existential quantifier is in some way removed from the range of the universal quantifier. That is, a resolution for those reflexive contradictions determined by C is:

**R**    Where $A(x,y)$ is a formula with just two distinct variables $x$ and $y$, both of which are free, and is such that $A(x,y) \supset (x \neq y)$ is a thesis, the incon-sistency arising from the assumption $(\exists x)(y)A(x,y)$ will be removed if, and only if, the value chosen for $x$ is removed from the range of the universal quantifier over $y$.

This is essentially Frege's conclusion. But it is important to see that in the form in which it is expressed here, it is a metaconclusion. By this is meant that it is a conclusion about *how* a resolution is to be achieved, but it is not a conclusion which determines the actual modification which has to be made in order to block the contradiction. Thus it determines the form of a solution, but not the manner of its expression. Hence the condition that the value chosen for $x$ has to be removed from the range of the universal quantifier is at this stage unavoidably vague. In particular, the word 'remove' should not be taken to imply that the value chosen to instantiate the existential quantifier is not formally available as an instantiation value for the universal quantifier and hence that a resolution is possible only in a many-sorted theory. What is intended, rather, is a weak sense of 'remove' which is satisfied if, say, some restrictive condition is put on the universal quantifier which the existential instantiation value fails to satisfy in particular cases. The restriction may occur explicitly in the quantifier, (e.g., as in many-sorted theories), or in the formula (e.g., as an antecedent condition), or in the inference rules (e.g., as in condi-tional universal instantiation); or it may be secured implicitly by some meta-condition (e.g., as in type theory). Thus a large variety of both formal and metadevices satisfy the metalinguistic requirement R.
    It is for this reason that there are many different competing solutions of the paradoxes. However, to the extent that differing solutions are genuine resolutions, they are not competing, for each succeeds because, and only because, it imposes an appropriate restriction on the quantification. But since a number of techniques do satisfy R, and since different ways of restricting the quantification have different consequences beyond the shared consequence of removing the inconsistency, many of which are related to the problem of plausibility, the question arises as to whether or not one technique is better than another; and this can only be answered in terms of the rationale which justifies the kind of quantifier restriction which is favoured. Thus, usually, a solution is a resolution together with a rationale.

In a sense, therefore, R tells us nothing since it neither points to a technique nor to a rationale. On the other hand, since C provides us with a purely formal criterion for the occurrence of reflexive contradictions, it seems reasonable to regard the problem of paradox removal as itself being a purely formal problem to be resolved only by whatever formal technique removes the contradiction, has no further consequences (such as ruling out too much), and is independent of the problem of plausibility. For given that a formula with free variables $x$ and $y$ entails that $x \neq y$, and given that the quantification over that formula commits us to the presupposition that $x = y$ in at least one instance, there is a straightforward formal inconsistency to be removed. Moreover, we know that no course other than a modification of the quantificational conditions is open to us,[5] however this may be done. Hence, it may seem, if our concern is strictly formal, then the best possible solution will be a minimal resolution, i.e., one which removes the contradiction and has no further consequences. To interpret R in this way would be to limit the number and kind of ways in which it can be satisfied and it would also avoid the need for an extraneous rationale. Whether or not there is such a minimal resolution will be taken up in Part II, Section 4.1.

R is of course limited in its application since it applies only to formulas which satisfy condition C. In the case of the general condition $C_1$ there is, correspondingly, a general form of R, though it is not easy to state in specific terms. Thus,

$R_1$    Where $QA$ is a formula satisfying condition $C_1$, the inconsistency arising from the assumption $QA$ will be removed if, and only if, the quantificational conditions are modified in such a way that the instantiation cases which are incompatible with the inequalities entailed by $A$ no longer arise.

There is an additional unsatisfactory vagueness here, beyond the vagueness present in R, due to the intuitive formulation of $C_1$. A more precise formulation will be given in Part II, Section 3.1.

**4 Summary**    Our main purpose in Part I has been to give a general formal characterization of reflexive contradictions in order to diagnose the cause and to prescribe a remedy. In the simple cases, the cause is given by C and the form of a remedy, though not a particular recipe, is given by R. More generally, given any sentence which can be expressed in the symbolism of quantification theory, it will be a reflexive contradiction in case it satisfies $C_1$; and the contradiction will be removed if $R_1$ can be satisfied. At this stage, however, these last two remarks are no more than intuitive claims which have yet to be made more precise.

Our interest in looking for a general characterization is motivated by a concern to understand how and why the simple paradoxes arise. By putting them in a wider context in which the meaning of the predicate constants plays no part we can see that the contradiction has nothing to do with the nature of classes, adjectives, barbers, catalogues, or falsity. Instead, it has to do with the compatibility of certain kinds of quantificational structure and standard identity conditions.[6] Formulas whose quantificational structures are incompatible with inequalities entailed by their matrices are demonstrable contra-

dictions since their negations are theses of pure quantification theory. But at this level there is no paradox, for there is no difference, from a formal point of view, between saying that $p$ & $\sim p$ is a contradiction and saying that $(\exists x)(y)$ $(f(y,x) \equiv \sim f(y,y))$ is a contradiction. The negations of each are theses, but quantification theory is consistent. So there is no inconsistency at this level. The inconsistency arises when, and only when, the negated thesis is assumed as a premiss; and exactly this is a feature of the familiar paradox arguments.

What is paradoxical about the paradoxes, therefore, is not that contradictions are demonstrable from the assumptions which are made, for the assumptions are themselves contradictory, but that there are informal, intuitive grounds for thinking that the assumptions are or ought to be true. Thus they seem to be inherently plausible. However, this feature does depend crucially on the meaning of the predicate constants involved and it is characteristic of only a small number of "paradoxical" sentences which satisfy C, or more generally $C_1$. Thus, it is plausibility not contradictoriness which depends on the nature of classes, adjectives, barbers, etc. But since the plausible assumptions *are* contradictions, the plausibility must be spurious (given that we are not going to give up quantification theory). For even though we have a remedy in terms of which the contradictions can be removed, the nature of the removal cannot be such as to supply us with consistent definitions of $R$, $C$, *Het*, etc. These are and will remain inconsistent concepts. In this case, what does the removal of the contradictions amount to and in what way is the apparent plausibility of the concepts spurious? An attempt to answer these questions is made in Part II, Section 4.

## NOTES

1. In most of what follows the discussion is presented in terms of $\alpha$, and $\alpha$-instances, rather than $\beta$. This is no loss since consequences of $\alpha$ are consequences of $\beta$.

2. Some of our friends (but not all) find this claim counterintuitive. We remain unmoved. But for those who think we have failed to give a credible example of an implausible paradox here, we suggest they try 'bigger than' or 'underneath' instead of 'to the right of'. In the case of 'bigger than' the postulated $\alpha$-instance affirms the existence of an ultimate smallest object; in the case of 'underneath', an ultimate uppermost object.

3. The *word Het* exists and the liar *sentence* exists, but what the contradictions show is that the word *Het* cannot have the sense intended for it and the liar sentence cannot yield the statement it is supposed to yield (cf. [5], pp. 108-109).

4. For convenience, $\epsilon(y,x)$ is now written $y \epsilon x$, but there is still no assumption that $\epsilon$ is characterized by set-theoretic postulates. In particular, no abstraction principle is presupposed. It is not necessary to use an abstraction principle to generate any of the class paradoxes.

5. Strictly, of course, since the inconsistency arises from an incompatibility between the quantificational structure of the formula and the identity conditions entailed by the matrix, there is a theoretical alternative to a modification of the quantificational conditions, namely a modification of standard identity criteria. But this does not seem to us to be a real option in the general case, though in the special case of the class paradox

this was of course Frege's way out—more precisely, Frege runs both options together by restricting the quantifiers in the identity criterion for classes.

6. In this connexion it is interesting to note that Thomson's theorem can be proved directly from the identity condition $x = y. \supset. f(y,x) \equiv f(y,y)$, since we have $(\exists y)(x = y) \supset (\exists y)(f(y,x) \equiv f(y,y))$, hence by MP and Gen, $(x)(\exists y)(f(y,x) \equiv f(y,y))$, i.e., $\sim(\exists x)(y) (f(y,x) \equiv \sim f(y,y))$. This demonstrates the essential dependence of the paradoxes on identity conditions.

## REFERENCES

[1] Goldstein, L., "Four alleged paradoxes in legal reasoning," *Cambridge Law Journal*, vol. 38, no. 2 (1979), pp. 373-391.

[2] Herzberger, H., "Paradoxes of grounding in semantics," *Journal of Philosophy*, vol. 67, no. 6 (1979), pp. 145-167.

[3] Martin, R. L., "On a puzzling classical validity," *Philosophical Review*, vol. 86 (1977), pp. 454-473.

[4] Ramsey, F. P., *The Foundations of Mathematics*, Routledge & Kegan Paul, London, 1931, pp. 20-21.

[5] Thomson, J. F., "On some paradoxes," pp. 104-119 in *Analytical Philosophy* (First Series), ed. R. J. Butler, Blackwell, London, 1962.

*University of Melbourne*
*Parkville, Victoria 3052*
*Australia*

*and*

*Princeton University*
*Princeton, New Jersey 08544*