

GRAHAM PRIEST

THE HOODED MAN

Received in revised version 14 May 2002

ABSTRACT. The Hooded Man Paradox of Eubulides concerns the apparent failure of the substitutivity of identicals in epistemic (and other intentional) contexts. This paper formulates a number of different versions of the paradox and shows how these may be solved using semantics for quantified epistemic logic. In particular, two semantics are given which invalidate substitution, even when rigid designators are involved.

KEY WORDS: epistemic logic, Eubulides, substitution of identicals

1. EUBULIDES THE PARADOXER

The most famous paradoxer of Antiquity is undoubtedly Zeno. His paradoxes, particularly those of motion, have exercised philosophers since he formulated them. But to my mind, the greatest paradoxer of Antiquity was not Zeno but the Megarian philosopher Eubulides. Eubulides is reputed to have formulated seven paradoxes, which Diogenes Laertius lists as: the Liar, the Disguised, the Electra, the Veiled Figure, the Sorites, the Horned One, and the Bald Head.¹ It would appear that some of these were variants of the others, and that there were basically four different paradoxes, which are as follows:²

1. *The Liar.* 'A man says that he is lying. Is what he says true or false?'
2. *The Hooded Man, the Unnoticed Man, the Electra.* 'You say you know your brother. But that man who came in just now with his head covered is your brother, and you do not know him.'
3. *The Bald Man, or the Heap.* 'Would you say that a man was bald if he had only two hairs? Yes. Would you . . . , etc. Then where do you draw the line?'
4. *The Horned Man.* 'What you have not lost you still have. But you have not lost horns. So you still have horns.'

Eubulides' arguments must have seemed like sophisms to many of his contemporaries, and made him an easy target for parody. Indeed, a contemporary Comic poet wrote:³

Eubulides the Eristic, who propounded his quibbles about horns and confounded the orators with falsely pretentious arguments, is gone with all the braggadocio of a Demosthenes.



But from the perspective of two and a half thousand years later, this low opinion is hardly justified.

The fourth of the above paradoxes is certainly little more than a sophism. It employs a device that is often used by barristers and other tricksters, and would now be classified as a *Fallacy of Many Questions*, of the kind ‘Have you stopped beating your wife?’. Literally, if you never had horns then you never lost them. Thus, the conditional ‘If you have not lost horns you (still) have them’ is false. The trick gets its bite from the conversational implicatures generated by the sorts of context in which one would normally talk of loss. The first and third paradoxes, the Liar and the Sorites are, by contrast, quite different. As no one familiar with contemporary philosophical logic needs to be told, these are of central importance to contemporary debates. Moreover, two and a half thousand years since Eubulides, there is still no consensus at all as to how to solve either of these paradoxes. This attests to their profundity. Compare the situation with that concerning Zeno’s paradoxes. Though philosophers may still argue about them, there has been, for at least a century, a general consensus concerning the solution to these paradoxes. This is why I said that, of Zeno and Eubulides, it is the latter who is the greater.

2. THE HOODED MAN PARADOX

What of Eubulides’ second paradox, the Hooded Man? Is this profound, like the first and third, or trivial, like the fourth? At first blush, it looks more like the fourth. But it is not. Though it may not have the depths of the first two, I think that it is a significant and hard paradox. It is the topic of this paper.

Let us start with a clean formulation. We suppose that a man walks into the room. The man is wearing a hood, and unbeknownst to you, it is your brother. Then the argument is simply:

This man is your brother.
You do not know this man.
<hr style="width: 100%; border: 0.5px solid black;"/>
You do not know your brother.

The premises are true, but the conclusion appears to be untrue. Yet the argument is an instance of the Substitutivity of Identity (SI):

$$a = b, \alpha(a) \vdash \alpha(b).$$

It is of the same form as: This man is your brother; this man has red hair; hence your brother has red hair. And this certainly seems to be valid.

name referring to the object in question. And the context does not change in this argument, so we can simply ignore the extra complexity created by the demonstrative, and take the person to be referred to by a name. Let us, therefore, christen the hooded man ‘Nescio’, where this is a rigid designator. (That is, a term whose denotation remains the same in all worlds, times, situations.) The argument then becomes:

Nescio is your brother.
You do not know that Nescio was born in Megara.
 You do not know that your brother was born in Megara.

Or to contrapose, simplifying again:

Nescio is your brother.
You know that your brother was born in Megara.
 You know that Nescio was born in Megara.

In what follows, when I refer to the Hooded Man argument, it is this argument to which I will be referring.

And with this version of the argument, we cannot avoid the problem, as we did the original, by saying that you *do* know that Nescio was born in Megara; you just don’t realise this. For you certainly do realise that your brother was born in Megara. Hence, the problem would then reappear with a different example of the same kind:

Nescio is your brother.
You realise that your brother was born in Megara.
 You realise that Nescio was born in Megara.

– And if you still doubt, consider the fact that the people of the 13th century did not know that water was H₂O. They did know that water was water. Or that one knows *a priori* that George Elliot was George Elliot, but one does not know *a priori* that George Elliot was Mary Anne Evans.⁶

3. SEMANTICS OF EPISTEMIC LOGICS

An advantage of formulating the Hooded Man argument in the way that I have done is that it allows us to bring to bear standard possible-world semantics, as they are deployed in an epistemic logic. Let us start our analysis of the argument with these.

We suppose that we have a first-order language with an identity predicate (but, for simplicity, no function symbols). In addition, the language has a one place operator, *K*, which can be read ‘You know that’, or ‘It is

known that', since the agent doing the knowing is playing no significant role here.⁷

A standard interpretation for the language⁸ is a structure $\langle D, \delta, W, R \rangle$. D is a domain of individuals. δ assigns every constant a world-invariant denotation in D . δ also assigns each predicate an appropriate extension at each world. So if P is an n -place predicate, $\delta(w, P)$ is a set of n -tuples of members of D . $\delta(w, =) = \{\langle d, d \rangle : d \in D\}$. Note that this particular extension is world-invariant. W is a set of worlds, and R is a binary accessibility relation on W . Intuitively, wRw' means that w' is a world where things are compatible with how they are known to be at world w . That is, everything known at w is true at w' . In standard epistemic logics, R is required to be reflexive, and maybe possess other properties too,⁹ but none of this will bear on matters here.¹⁰

We now specify what it is for a sentence, α , to be true at a world ($w \Vdash \alpha$), by the familiar recursive clauses. For all $w \in W$:

$$w \Vdash Pt_1 \dots t_n \text{ iff } \langle \delta(t_1), \dots, \delta(t_n) \rangle \in \delta(w, P),$$

$$w \Vdash \neg\alpha \text{ iff } w \not\Vdash \alpha,$$

$$w \Vdash \alpha \wedge \beta \text{ iff } w \Vdash \alpha \text{ and } w \Vdash \beta,$$

$$w \Vdash \alpha \vee \beta \text{ iff } w \Vdash \alpha \text{ or } w \Vdash \beta,$$

$$w \Vdash \alpha \rightarrow \beta \text{ iff for all } w \in W \text{ such that } w \Vdash \alpha, w \Vdash \beta,$$

$$w \Vdash K\alpha \text{ iff for all } w' \text{ such that } wRw', w' \Vdash \alpha,$$

$$w \Vdash \exists x\alpha \text{ iff for some } c \in C, w \Vdash \alpha_c^x,$$

$$w \Vdash \forall x\alpha \text{ iff for all } c \in C, w \Vdash \alpha_c^x.$$

Here, C is the set of constants, and α_c^x is α with all free occurrences of 'x' replaced by 'c'. I assume, at this point, that the language has been augmented, if necessary, so that each member of the domain has a name. This is not essential, but merely relieves us of the necessity of talking in terms of satisfaction. Note that the conditional, \rightarrow , is the strict conditional of *S5*. Validity, \models , is defined in terms of truth preservation at all worlds of all interpretations.

These semantics have their problems, and it may well be that more sophisticated semantics are needed for various reasons, for example, to give a decent account of the conditional.¹¹ But they will do us to start with.

4. THE EASY VERSION OF THE ARGUMENT

There is still one more preliminary issue that needs to be addressed. This is how to understand the noun-phrase ‘your bother’. This is a demonstrative, but it is not a simple demonstrative, since it packs in the information that the thing referred to has a certain property – that of being a brother. Since the context is not changing, we can ignore the demonstrative aspect of it, but the question is whether the phrase should be understood as functioning like a name or like a description.

Suppose, first, that it is understood as a description, ‘the thing which is your brother’, $\iota x Bx$.¹² In this case, the argument is invalid. It is well known that SI may fail in modal contexts when the term substituted is a description.¹³ For example, it is necessarily the case that $9 = 9$; but it is not necessarily the case that the number of planets = 9. Though ‘the number of the planets’ refers to 9 in this world, in other worlds, it may refer to a different number.

To show how this works in detail for the Hooded Man argument, the language must be augmented with a description operator, ι , with the usual syntax. In the semantics, descriptions are given world-dependent denotations by the following recursive clause:

$$\delta(w, \iota x \alpha) = d \text{ if } d \text{ is the unique } d \in D \text{ such that for some name of } d, c, \\ w \Vdash \alpha_c^x$$

and, let us suppose, some fixed but arbitrary denotation otherwise. (This is hardly an adequate semantics to handle non-denotation. But again, the complexities of non-denotation are not relevant to the matter at hand, so we can keep it simple.) If c is a constant, let us define $\delta(w, c)$ as $\delta(c)$; the truth conditions of atomic sentences may then be given uniformly as:

$$w \Vdash Pt_1 \dots t_n \text{ iff } \langle \delta(w, t_1), \dots, \delta(w, t_n) \rangle \in \delta(w, P).$$

We can now represent this version of the Hooded Man argument as follows:

$$\begin{array}{l} n = \iota x Bx \\ \hline KM \iota x Bx \\ \hline KMn \end{array}$$

The argument is formally invalid. The following is a counter-model.

$$W = \{w, w'\},$$

$$wRw',$$

$$D = \{0, 1\},$$

$$\delta(n) = 0,$$

$$\delta(w, B) = \delta(w, M) = \{0\},$$

$$\delta(w', B) = \delta(w', M) = \{1\}.$$

As is easy to check, $\delta(w, \iota x Bx) = 0$, and $\delta(w', \iota x Bx) = 1$. In particular, then $w \Vdash n = \iota x Bx$. Moreover, $M \iota x Bx$ is true at w and w' . Hence $w \Vdash KM \iota x Bx$. But Mn is false at w' . Hence, $w \not\Vdash KMn$.

In this version of the argument, the description has narrow scope. One may well ask what happens to the argument if the description is interpreted as having wide scope. I will return to this question in due course.

5. THE HARD VERSION OF THE ARGUMENT

So much for interpreting the phrase ‘your brother’ as a description. Let us now turn to the other possibility mentioned, that is, interpreting the phrase as picking out an individual rigidly (referring to the same thing in every world). In this case, we can treat it as a name. Indeed, suppose that your brother’s name is ‘Cain’. We can simply use this name. The argument is now:

Nescio is Cain.
 You know that Cain was born in Megara.
 —————
 You know that Nescio was born in Megara.

This is of the form:

$$\begin{array}{l} n = c \\ \frac{KM c}{KM n} \end{array}$$

And it is not hard to check that in the preceding semantics this argument is formally valid.

Yet this seems wrong. You *do* know that Cain was born in Megara. You *don't* know that Nescio was. Hence, despite the formal semantics, this form of SI does seem invalid. How so?¹⁴

6. THE PUZZLE ABOUT PIERRE

Before we look at answers to this question, it is worth noting that the view that examples of the kind with which we are dealing demonstrate the failure of SI in epistemic contexts has been challenged by Kripke (1979). He argues that contradictions of the kind in question arise even without SI. It is therefore wrong to point the finger of suspicion at it. His well-known example concerns a native French speaker, Pierre, who expresses one of his beliefs by saying 'Londres est jolie'. He then learns English, and comes to express one of his beliefs by saying 'London is not pretty' – without revoking any former dispositions concerning his assertions in French: he simply does not realise that 'Londres' and 'London' refer to the same place. He would appear to believe that London is both pretty and not pretty. He may even vehemently reject the claim 'London is pretty', in which case he would seem both to believe and not to believe that London is pretty.

In fact, the detour through French is unnecessary, as Kripke points out. The issue would be just the same as that which arises for a person, say Pierre, who sincerely asserts: 'George Elliot is a man' and 'Mary Anne Evans is not a man' – or, even stronger, denies 'Mary Anne Evans is a man' – unaware that they are the same person.

Let us concentrate on the monolingual version. In the given situation, it is virtually irresistible to hold that:

Pierre believes George Elliot to be a man

and that:

Pierre believes Mary Anne Evans not to be a man

– or, in the stronger case, that:

Pierre does not believe Mary Anne Evans to be a man

We have this from the horse's mouth; and though this sort of evidence may be overridden in some cases (e.g., when speakers do not properly understand the words they use) we can set up the situation in such a way that cases of this kind are explicitly ruled out. Contradiction arises here when, and only when, we add the further premise that Elliot is Evans, to conclude that Pierre believes Elliot not to be a man – or, in the stronger case, does not believe Elliot to be a man. SI is essentially involved in these contradictions.

The contradictions that Kripke points to, by contrast, concern how Pierre's beliefs may be reported in *paraphrase*. If we paraphrase Pierre's beliefs about Evans by using the name 'Elliot' we have similar contradictions. Now we often paraphrase people's views in reporting them. Suppose that you tell me that it was the author of the Sherlock Holmes stories who was the Ripper. It would normally be fair for me to report your belief to a third party by saying that you think that Doyle was the Ripper. If, however, you thought also that Doyle stole the Holmes stories, this would not be fair paraphrase.

Similarly, and closer to hand, suppose that it is common knowledge in a group (which includes Pierre) that Elliot is Evans, and that Pierre believes Elliot to be a man (perhaps believing that Elliot was, in fact, a very successful transvestite), then it would be quite legitimate to report his belief by saying that he believed Evans to be a man. But in the sort of case in question, where Pierre does not know that Elliot and Evans are one, it would be quite misleading to paraphrase his belief that Evans is a woman by saying that he believes Elliot to be a woman.

The constraints on legitimate paraphrase, and, in particular, the role that background knowledge plays in the matter, are, I suspect, complex. But this is not the place to go into them. It is clear that Kripke's problem arises because of a violation of these constraints. The contradictions that we are concerned with do not depend in any way on paraphrase; and SI is central to them.

7. FREGE AND SI

That SI fails in epistemic contexts, even when names are involved, is, of course, a well-known view. It was Frege's.¹⁵ According to Frege's view, the Hooded Man inference fails because in the sentence 'You know that Cain was born in Megara', 'Cain was born in Megara' refers not to its standard reference (which is, for Frege, a truth value), but to its sense, the thought (proposition) that Cain was born in Megara. And 'Cain' refers, not to its standard referent, the person, but to its standard sense, something like a conception of that person (an individual concept). Similarly, in the epistemic context, 'Nescio' refers not to the person but to a conception of a person. And even if Nescio and Cain are the same, the two conceptions of the person are not. Hence, we cannot substitute the one for the other. (In a sense then, the failure of SI is merely syntactic, since we are not dealing with co-referring expressions.)

Unfortunately, Frege's account faces difficult problems.¹⁶ Consider the inference:

You know that Cain was born in Megara.
Cain has red hair.

There is someone with red hair whom you know
to have been born in Megara.

This would certainly seem to be valid, but it is not for Frege. Even if the premises are true, the conclusion:

(1) $\exists x (x \text{ has red hair} \wedge \text{you know } x \text{ was born in Megara})$.

is false. To make the sentence true, the first x has to be the person; the second x has to be an individual concept. And no person is an individual concept.

There are ways one might try to get around this problem. For example, one may bite the bullet and agree that (1) is really false. The truth that the conclusion is meant to express is:

$\exists x \exists y (x \text{ has red hair} \wedge y \text{ is an individual concept of } x \wedge \text{you know that } y \text{ was born in Megara})$.

But if this is the conclusion of the argument, it would still seem to be invalid. For the conclusion now entails the existence of individual concepts; but the premises certainly don't appear to do this.¹⁷

A somewhat different objection to Frege's account is as follows. Suppose that Arthur does not know the identity of Jack the Ripper. That is:

$\neg \exists x \text{ Arthur knows that } x \text{ is Jack the Ripper}$.

Or to put it in kosher Fregean terms:

$\neg \exists x \exists y (y \text{ is an individual concept for } x \wedge \text{Arthur knows that } y \text{ is Jack the Ripper})$.

Now, the individual concept *Jack the Ripper* is a concept for Jack the Ripper, and Arthur certainly knows that Jack the Ripper is Jack the Ripper. Hence:

$\exists x \exists y (y \text{ is an individual concept for } x \wedge \text{Arthur knows that } x \text{ is Jack the Ripper})$.

So Arthur does know the identity of Jack the Ripper.

Maybe there are ways to try to get around these difficulties.¹⁸ But the discussion suffices to show that there are enough problems with Frege's account¹⁹ to make it worth considering whether there are other, and simpler, semantic analyses of intentional contexts. There are; and to these we now turn.

8. A MORE SENSITIVE SEMANTICS FOR IDENTITY

I will give two. Both, it would seem, have a certain naturalness, and I am not sure which is preferable.

For the first semantics, come back to the Hooded Man, and consider the situation concerning Cain and Nescio as you know it. These two might have the same identity; they might not. That is, there are worlds/situations compatible with all that you know in which they do have the same identity, and worlds/situations in which they do not. In particular, then, it is possible for objects to have different identities in different worlds/situations. For any object, then, there is a function that maps it to its identity at each world. Indeed, more simply, we can just think of an object as a function that maps each world to an identity.

Formally, the semantics look like this.²⁰ An interpretation is a structure $\langle D, I, \delta, W, R \rangle$. W is a set of worlds, and R is a binary accessibility relation on W , as before. I is a set of things that we may think of as identities. D is a collection of functions from worlds to I ; so that if $f \in D$, $f(w)$ is the identity of f at w . δ assigns every constant a world-invariant denotation in D . (And we assume, as before, that every member of D has a name in the language.) δ also assigns each predicate an appropriate extension at each world. But the extensions now are subsets of (n -tuples of) I , not D .²¹ In particular, $\delta(w, =) = \{\langle i, i \rangle : i \in I\}$. (This extension is still world-invariant.)

Truth values are assigned to atomic formulas by the clause:

$$w \Vdash Pt_1 \dots t_n \text{ iff } \langle \delta(t_1)(w), \dots, \delta(t_n)(w) \rangle \in \delta(w, P).$$

The recursive truth conditions are the same as before. Note, in particular, that the domain of quantification is still D , not I . Validity is also defined the same way.

It is tempting to think of the values of a function, f , in D as the *parts* of the object f at each world. In the same way, if this were a temporal logic, it would be natural to think of f as an object comprising temporal parts, and as the members of I as the temporal parts. And one can certainly

conceptualise things in this way. I think that this is the wrong way to think about things, at least in the epistemic case, however. This way of looking at matters takes the parts to be metaphysically primary, and the object to be the sum of its parts. I think that it is preferable to take the members of D to be metaphysically primary. Their values at each world are their identities there. At each world an object has an identity, just as much as it has a length, a colour, and so on. (All the worlds are stages, and all the people merely players.) One argument for this is as follows. If there were parts that were metaphysically primary, there would appear to be no reason why every function from worlds to parts should not constitute an individual. (There are no privileged linkages between world parts.) But if this is the case, as is easy to check, the following would be a valid inference: $K\exists x\alpha \vdash \exists xK\alpha$. But this is certainly not valid. I can know that there are spies without knowing of any person that they are a spy.

Whatever one makes of these issues, since the truth conditions for connectives and quantifiers have not changed, the propositional/quantificational logic of these semantics is still the same. In particular, then, all the standard quantificational rules, such as Existential Generalisation – which fails on Frege's account – are valid.

The novel feature of these semantics shows up in the behaviour of identity. In particular, $\not\models \forall x\forall y(x = y \supset Kx = y)$. As a counter-model for this, take the structure where:

$$W = \{w, w'\},$$

$$wRw',$$

$$I = \{0, 1\},$$

$$D = \{f, g\}, \text{ where } f(w) = g(w) = f(w') = 0 \text{ and } g(w') = 1,$$

$$\delta(c) = f, \delta(n) = g.$$

Then since $f(w) = g(w)$, $w \Vdash c = n$. But since $f(w') \neq g(w')$, $c = n$ is false at w' , and so $w \not\models Kc = n$.

The solution to the hard Hooded Man problem is now simple. The inference in question ($n = c, KMc \vdash KMn$) is invalid. To see this, take the same interpretation where, in addition:

$$\delta(w, M) = \delta(w', M) = \{0\}.$$

For future reference, call this interpretation \mathcal{I}_1 . As before, $w \Vdash n = c$. And since $f(w) = 0 \in \delta(w, M)$, Mc is true at w . In the same way, it is

true at w' . Hence, $w \Vdash KMc$. But since $g(w') = 1 \notin \delta(w', M)$, Mn is not true at w' , and so $w \not\Vdash KMn$. Thus is the hard version of the Hooded Man solved.

It should be observed that in these semantics SI does hold provided that substitution occurs in non-epistemic contexts. The truth conditions for atomic sentences ensures this for atomic contexts, and an induction over the connectives and quantifiers (other than K), does the rest.²²

9. IMPOSSIBLE WORLDS

I now turn to a second semantics which invalidates SI. For this, we turn for inspiration to the semantics of relevant logics. In standard modal propositional logics, $q \rightarrow q$ is a logical truth, true at all worlds. Hence, $p \rightarrow (q \rightarrow q)$ is also a logical truth. (\rightarrow is the strict conditional here.) That is, there are “fallacies of relevance”: logical truths of the form $\alpha \rightarrow \beta$ where α and β share no propositional parameter. In relevant logic, there are no such logical truths. The major device in the world-semantics for relevant logics that delivers this effect is the employment of a new kind of world. We still have the same worlds as before, the possible worlds, or the *normal* worlds as they are usually called in this context. But we now add a bunch of *non-normal* worlds. These are thought of as (logically) impossible worlds. The idea that there can be physically impossible worlds, that is, worlds where the laws of physics may be different, is a fairly standard one. Such worlds are still logically possible. But just as there can be worlds where the laws of physics may be different, so there are worlds where the laws of logic may be different. Intuitively, after all, we reason about such worlds when we consider alternative logics. Thus, a classical logician believes that the law of excluded middle is valid. But they know well that if intuitionist logic were correct, this law would fail, though the law of non-contradiction would not. They therefore seem to be quite capable of considering logically impossible situations, and making discriminations about what happens within them. And given such non-normal worlds, p may hold at one of them where $q \rightarrow q$ fails. Hence, $p \rightarrow (q \rightarrow q)$ is not a logical truth. ($\alpha \rightarrow \beta$ holds at a normal world if every world (normal or non-normal) where α holds β holds.²³)

So much for the basic idea. The next question is how, as a matter of technique, one arranges for $q \rightarrow q$ and its like to fail at a non-normal world. There are, in fact, a number of different techniques that can be deployed here.²⁴ A prominent one is using a ternary relation to give the truth conditions of \rightarrow . However, the simplest is just to assign conditionals *arbitrary* truth values at non-normal worlds. After all, conditionals of

entailment strength express laws of logic, and if logic is allowed to vary, such conditionals may behave in any way.

Let us now return to the Hooded Man argument. The invocation of impossible worlds has no effect on this. The conditional does not feature in the argument, so messing around with its behaviour is irrelevant. However, we can adapt the techniques just reviewed. Some of the worlds in the semantics of an epistemic logic represent the way the world could be, as far as is known. Now, arguably, that you know α does not entail that you know *anything* else. For example, suppose that you know that α . It does not follow that you know $\alpha \vee \beta$. If β , for example, employs concepts that are unknown to you (in the way that the concept of a micro-processor was unknown to a medieval monk), then you do not believe $\alpha \vee \beta$. *A fortiori*, you do not know it. Similarly, suppose that you know that α . It does not follow that you know that $\neg\neg\alpha$. You may believe the law of double negation to be invalid. Hence you may not believe $\neg\neg\alpha$. *A fortiori*, you do not know it. Similar objections can be brought against similar examples. And if they stand, then there may be worlds where there is no essential connection between the holding of different sentences; any sentence may hold or fail, independently of any other. Technically, then, at these worlds we may assign arbitrary truth values to *all* sentences, not just to conditionals.

Making these ideas precise: an interpretation for the language is now a structure $\langle D, \delta, W, W^+, R \rangle$. Everything is exactly the same as in the basic epistemic semantics of section 3 – and, note, identity still has its usual semantics – except for two things. First, $W^+ \supseteq W$. The members of $W^+ - W$ are those worlds that represent states of knowledge other than those represented by standard possible worlds. We might call these *really non-normal* worlds.²⁵ Secondly, for every $w \in W^+ - W$, δ assigns every formula, α , a truth value, $\delta(w, \alpha)$ (T or F). Here, note, the assumption that every element of the domain has a name is essential. This is *not* a shorthand for talking in terms of satisfaction. Quantification is therefore, in effect, substitutional.

The truth conditions at normal worlds are as before. But if $w \in W^+ - W$:

$$w \Vdash \alpha \text{ iff } \delta(w, \alpha) = T.$$

Validity is still defined in terms of truth preservation over the worlds in W , not W^+ .

The really non-normal worlds have no effect on inferences involving non-epistemic notions. The non-epistemic fragment of the logic is still, therefore, classical first-order logic. The worlds do have an effect on inferences concerning K , however. In particular, SI fails in epistemic contexts.²⁶

The following is a counter-model to the hard version of the Hooded Man argument. $W = \{w\}$, $W^+ = \{w, w'\}$, wRw' , $D = \{0\}$; δ is such that $\delta(n) = \delta(c) = 0$, $\delta(w', Mc) = T$, and $\delta(w', Mn) = F$. It is almost trivial to check that $w \Vdash n = c$, $w \Vdash KMc$, but $w \not\Vdash KMn$. For future reference, let us call this interpretation \mathcal{I}_2 . Note that, as in the semantics of the previous section, SI holds provided that we are not substituting into the scope of a K .

One might well worry that this solution is too cheap. It destroys all inferences concerning knowledge;²⁷ and knowledge is not that anarchic. For the reasons I have given, I think it is. But note, in this context, the following. First, all intensional verbs pose problems of substitutivity of exactly the same kind. One might believe that Nescio was born in Megara without believing that Cain was. Or fear, or hope, or doubt, etc. A *general* solution to the problem must therefore apply to all such verbs. And most of these are very anarchic. I have already noted that one can believe anything whilst, at the same time not believe pretty much anything else. Similarly, one can fear, quite rationally, that $\alpha \wedge \beta$ without fearing that α or fearing that β (e.g., I might not fear Jack's going to the party, nor fear Jill's going to the party, yet still fear that Jack *and* Jill go to the party – because there will be a terrible argument). And with irrational fears, all bets are off. So the semantics provides a *unified* solution to all these problems of substitutivity.

Next, if some of these intensional notions do have more inferential structure, this can be recaptured by adding constraints on how δ behaves. Thus, suppose that it is thought we are dealing with a logically perfect agent, whose knowledge is closed under logical consequence, we simply require the set of truths at each member of $W^+ - W$ to be closed under logical consequence; or if we think that the inference $Ka = b, K\alpha(a) \vdash K\alpha(b)$ is valid, we simply close the things assigned true at every such world under the identities that are true there.²⁸ And all such constraints can be imposed without threatening the solution to the problem. All that the solution requires is that it be possible to assign Mc and Mn different truth values at a world, even though 'c' and 'n' actually refer to the same person. And the legitimacy of this derives from the fact that the world realises the way the agent represents the world to be.²⁹

10. THE *De Re* ARGUMENT

We now have two semantics which deliver the failure of SI in epistemic contexts, even when rigid designators are involved.³⁰ Which, if either, is the better one, I leave the reader to decide. We are still not finished with the Hooded Man argument.

There is a distinction, dating back to Medieval Logic, that is standardly drawn between two different understandings of a statement of the form 'It is known that Cain was born in Megara'. On the first understanding, *de dicto*, this expresses a property of a proposition, or some other kind of truth-bearer, such as a sentence; in this case, the proposition (or sentence) *Cain was born in Megara*. The epistemic sentences we have been concerned with so far are, in fact, all of this kind.

On the second understanding, *de re*, the sentence is taken to express a predication of the object of the belief, in this case, Cain. The *de re* interpretation might be expressed more perspicuously as: Cain is such that he is known to have been born in Megara. It is usually claimed that SI holds for *de re* interpretations. Indeed, it is often taken as a criterion for being *de re*. Thus, there would appear to be another version of the paradoxical argument in the wings, which is as follows:

Nescio is Cain.
 Cain is such that you know that he was born in Megara.
 —————
 Nescio is such that you know that he was born in Megara.

Is the conclusion of this argument unacceptable though? Perhaps not. Nescio, *that very person*, is such that you know him to have been born in Megara. You just don't realise this.

But things are not that simple. Suppose that SI works in *de re* contexts. Then it is indeed true that Nescio is a person, viz., Cain, such that you know him to have been born in Megara. But it would appear equally to be the case that Cain is a person, viz. Nescio, such that you do not know him to have been born in Megara, since we have:

Nescio is Cain.
 Nescio is such that you do not know him to have been born in Megara.
 —————
 Cain is such that you do not know him to have been born in Megara.

Let us call this the *counter-argument*. It would seem to be just as good. And if so, there is a person (Nescio, i.e., Cain) such that you both know and do not know him to have been born in Megara. We still appear to have a contradiction on our hands. What is to be said of this?

One possible solution to the problem is to insist that the second premise of the counter-argument, that Nescio is such that you do not know him to have been born in Megara, is just false. He *is* such that you know him to have been born in Megara; you just do not realise this fact. You may not know, *de dicto*, that Nescio was born in Megara. But what you know about

Nescio *de re* is not open to introspection, simply because you may not recognise him under certain descriptions.

This is certainly a possible solution, but it has its problems. We have granted the *de dicto* claim that you do not know that Nescio was born in Megara. Moreover, ‘Nescio’ here is a rigid designator. It refers to that very object, independently of how it is picked out in a particular world (unlike the way a description may refer). There may even be a causal (indeed perceptual) baptism of Nescio with this name. It would seem to follow that the epistemic state is also *de re*. That is, *de dicto* + (rigid designation) + perceptual contact entails *de re*.³¹

Perhaps there are replies to this objection. The notion of *de re* knowledge is, after all, slippery enough. But is there another possible solution? There is. To see what it is, first consider the sentence:

Cain is such that you know that he was born in Megara.

What is its logical form? The way to represent the sentence which sticks most closely to its surface form is obtained by employing λ -abstraction, so that it may be represented as $\lambda x(KMx)c$.³² But we can avoid introducing this new machinery. The point of a *de re* claim is that it is a claim about the object itself, independently of how it is referred to. And since reference to objects themselves is carried by quantifiers, we can capture the content of the claim by:

$\exists x(x = \text{Cain} \wedge K x \text{ was born in Megara})$

Hence, the *de re* inference is of the form:

$$\frac{n = c \quad \exists x(x = c \wedge KMx)}{\exists x(x = n \wedge KMx)}$$

This argument, involving substitution, as it does, only in non-epistemic contexts is valid in both the previous semantics. Given that the premises are true, we therefore accept the conclusion: the hooded man, Nescio, is such that he is known to have been born in Megara. (Though, as referred to by the name ‘Nescio’, you may not realise this.)

What of the counter-argument? The logical form of the argument is:

$$\frac{n = c \quad \exists x(x = n \wedge \neg KMx)}{\exists x(x = c \wedge \neg KMx)}$$

and it, too, is valid. Hence, if the premises are true, so is the conclusion: Cain is such that you do not know him to have been born in Megara. (Though, as referred to by the name ‘Cain’, you may not realise this.)

Thus Nescio (that is, Cain), is such that he both is and is not known to have been born in Megara. This may *sound* like a contradiction, but it is not. It is of the form:

$$(3) \exists x(x = n \wedge KMx) \wedge \exists x(x = n \wedge \neg KMx).$$

Of course, the x in question is n , and $\neg KMn$; but any attempt to obtain KMn , and thus an explicit contradiction, from the first conjunct, falls foul of the failure of SI in epistemic contexts. Indeed, both interpretations \mathcal{I}_1 and \mathcal{I}_2 of previous sections make (3) true, since they make $c = n \wedge KMc$ and $n = n \wedge \neg KMn$ true at world w . The result follows by existential generalisation.³³ Hence, the *de re* problem is also solved by both semantics.

Let us return, finally, to the interpretation of the argument which employs a description of wide scope. In this case, the sentence ‘You know your brother to have been born in Megara’, can be represented as $\exists x(x = \iota yBy \wedge KMx)$ (or as a similar thing employing λ -terms). As such, it is essentially of the same form as the *de re* statement using a name. In fact, in a *de re* context, how an object is specified is irrelevant, since we are talking about the satisfaction of a condition by an object itself. Thus, the analysis of the argument in this case is essentially the same. The object both is and is not known to have been born in Megara, but this is not a contradiction, and the failure of substitutivity in epistemic contexts prevents it collapsing into one.

11. CONCLUSION

There was a sophisticated discussion of epistemic operators, and of Eubulides’ paradox, in Medieval logic,³⁴ though, like so much in logic, this disappeared with the rise of Modern philosophy. The issue was put back on the map by Frege, though it was quickly taken off again by what came to be the dominant extensionalism of the *Tractatus* and, later, of Quine.³⁵ The topic of intensional discourse was again replaced on the map by the development of possible-world semantics in the 1960s and 1970s – though the discussion of identity and epistemic contexts has still received much less attention than the problem of substitutivity in modal contexts. I have argued that an adequate analysis of the semantics of epistemic (and similar intentional) notions requires a robust failure of SI, and suggested two semantics which deliver this. Both invoke machinery that goes beyond that

which is standardly employed in modal logic. What Eubulides would have made of all this, I have no idea. But I guess that he could not have failed to be pleased by the fact that the Hooded Man paradox – along with the Liar and the Sorites – are still being discussed by logicians over two thousand years later.³⁶

NOTES

¹ Hicks (1925): II, 108. Naturally, one can dispute whether Eubulides really did invent these paradoxes. For example, some have attributed the Disguised (the Hooded Man) to Euclides, the founder of the Megarian school.

² See Kneale and Kneale (1962), p. 114, who cite the classical sources.

³ Hicks, *loc cit.*

⁴ See Priest (2000).

⁵ Hintikka (1962), p. 132, claims that ‘you know who *a* is’ is to be understood as: $\exists x Kx = a$. This is dubious. Not only does it not make *knowing who* context-independent, but it requires the failure of existential instantiation, even if *a* is a rigid designator, since, presumably, it is always true that $Ka = a$. But even this understanding cashes out *knowing who* in terms of *knowing that*. The most careful analysis of *knowing who* of which I am aware is provided by Boër and Lycan (1986). They recognise the context-dependence of *knowing who*, taking the relevant contexts to be certain speaker-purposes. They argue that *knowing who* is a certain kind of *knowing that*, and provide a sensitive analysis of the kind of *knowing that* that it is. The details need not concern us here. It should be noted that their analysis is greatly complicated by their employment of a paratactic analysis of *knowing that*, rather than the much more straightforward one employed in what follows.

⁶ Salmon (1982) has an account of intensional contexts which allows substitutivity universally in such contexts. According to him, knowledge – and similar intensional predicates – are, *stricto sensu*, ternary relationships between an agent, a proposition, and a guise (or Fregean sense). Knowledge *simpliciter* is always knowledge relative to *some* guise. Hence, he thinks that you do realise that Nescio was born in Megara (since you know it relative to Nescio’s guise *your brother*). You do not know it relative to the guise *man who has just entered the room*. Moreover, Salmon claims, the proposition that George Elliot is George Elliot is *a priori*. Hence (*contra* Kripke), so is the proposition that George Elliot is Mary Anne Evans. It seems to me that someone who accepts this has lost contact with the work that a priority needs to do. There is no way that that particular fact could be reasoned out without empirical knowledge. One might be tempted to say that there is some guise under which it could not be reasoned out without empirical knowledge; but for Salmon, a priority is a property of propositions; it is not relative to a guise (p. 133). Nor can one say that it is the truth of the sentence ‘George Elliot is Mary Anne Evans’ that cannot be reasoned out without empirical knowledge: the truth of no sentence can be reasoned out without empirical knowledge (about meanings).

⁷ More generally, one could take *K* to be an operator requiring an agent and a sentence: $xK\alpha$ (*x* knows that α). Semantically, the interpretation of such an operator would be a family of binary relations, R_d , indexed by members of the domain of quantification, *d*.

⁸ See, e.g., Fitting and Mendelsohn (1989), Ch. 4.

⁹ See Priest (2001), 3.6.

¹⁰ These semantics are constant-domain. There are, of course, other world-semantics in which the domains may vary. I happen to think that constant-domain semantics are correct, and that the effect of variable domains should be obtained with an appropriate existence predicate. However, this is not the place to go into all this.

¹¹ A problem that has nothing to do with the conditional is the problem of “logical omniscience”: if $\alpha \models \beta$ then $K\alpha \models K\beta$. This seems far too strong. Even if you know α , if β is a very complicated consequence of α , you may not believe β , or *a fortiori*, know it.

¹² Since brothers may not be unique, this should be an indefinite description. But the extra complexity does not affect the situation, so we keep things simple.

¹³ See, e.g., Fitting and Mendelsohn (1998), esp. Chs. 9, 12.

¹⁴ One might maintain that all names are really covert descriptions, and thus reduce this version of the argument to the previous one. But such a move is well known to face grave difficulties. See Kripke (1980). It is particularly hard to suppose that demonstratives are covert descriptions, since, for these, uptake of reference may be secured with no linguistic intermediary.

¹⁵ See Frege (1952).

¹⁶ Similar problems beset a paratactic account of *knowing that* (of the kind developed by Boër and Lycan (1986)) and any other account which makes it impossible for a variable to bind inside and outside an epistemic context simultaneously in the natural way.

¹⁷ A similar argument is often employed in connection with plural predication and quantification. There are some sentences that contain plural predicates and quantifiers, and which cannot be cashed out in terms of standard first-order quantifiers. A notorious example is:

(2) There are some critics who admire only each other.

Some writers have suggested that this sentence is actually a covert second-order sentence, quantifying over a set of critics. Such a suggestion would certainly appear to be incorrect. (2) appears to entail the existence of critics, but not of sets – which the second-order sentence does. See, e.g., Yi (1999), esp. p. 165f.

¹⁸ For example, Kaplan (1971) tries to get around the problem about Arthur, by restricting the quantifier over senses to those that are “vivid”, where vividness is ‘intended to go to the purely internal aspects of individuation’ (p. 135). Unfortunately, no very precise characterisation of vividness emerges.

¹⁹ There are also, of course, other objections to Frege’s theory. The whole idea that proper names have a semantically significant sense has been attacked by Kripke (1980).

²⁰ After writing the first draft of this paper, I discovered that essentially these semantics for contingent identity systems are given by Parks (1974). Other systems for contingent identity can be found in Hughes and Cresswell (1968), p. 198f., Bressan (1972), and Gupta (1980). For a discussion of the first two of these, see Parks and Smith (1974), and Parks (1976), respectively. None of the above is concerned with epistemic contexts.

²¹ It is tempting to think of identities as Fregean senses, but this would not be right. If anything, it is the members of D that are more like senses, since they determine behavior across worlds. This is essentially how members of D are, in fact, interpreted by Bressan (1972), Gupta (1980), and, in a similar semantics, Hintikka (1971). It would also be a mistake to interpret members of D , in the present semantics, as senses, however. They are simply the objects themselves.

²² There is some resemblance between the above semantics and Lewis' counterpart theory (1968) – page references to the reprint. Thus, it might be thought that the functions in D are simply the counterpart relations between members of I . Specifically, we might think that x is a counterpart of y iff:

$$\exists f \in D \exists w_1, w_2 \in W f(w_1) = x \wedge f(w_2) = y.$$

This, however, is not the case. First, there are questions of interpretation. In counterpart theory it is the members of I that are the genuine objects, and so constitute the domain of quantification; in the above semantics it is the members of D that are the genuine objects. Secondly, there are differences between the properties of the above relation and a counterpart relation. The above relation is clearly symmetric, but a counterpart relation need not be (p. 28f.). Third, these differences affect the resulting logic. For example, in the above semantics, if we require R to be a universal relation, the resulting propositional modal logic is $S5$. It is not in counterpart theory. Specifically, the universal closure of $\alpha \supset \Box \Diamond \alpha$ fails since the counterpart relation is not symmetric (p. 36). Moreover, various quantification principles hold in the above semantics that fail in counterpart theory. For example, $\exists x \Box \alpha \supset \Box \exists x \alpha$ holds in the above semantics, but fails in counterpart theory (p. 36).

²³ Validity, however, is defined in terms of truth-preservation just a *normal* worlds. When we come to assess the validity of an argument, the logic in question must be the right one!

²⁴ For a full account of the matter, see Priest (2001), Chs. 8–10.

²⁵ More generally, an interpretation might have non-normal worlds of the more orthodox kind as well (so as to produce a relevant logic). If it has, then a conditional, $\alpha \rightarrow \beta$, should be true at a normal world, w , if the move from α to β is truth-preserving at normal and non-normal worlds – not at the really non-normal worlds as well. However, these matters have no effect on the problem at hand, so we may keep things simple.

²⁶ For good measure, the semantics also solve the problem of logical omniscience. Suppose that $\alpha \models \beta$. Consider an interpretation where $W = \{w\}$, $W^+ = \{w, w'\}$, wRw' , and δ is such that $\delta(w', \alpha) = 1$ but $\delta(w', \beta) = 0$. Then $w \Vdash K\alpha$, but $w \not\Vdash K\beta$. Hence $K\alpha \not\models K\beta$. After writing this part of the paper, I discovered that these semantics (with some minor and inessential modifications) were given by Rantala (1982). He proposes them to solve the problem of logical omniscience. The language he uses does not contain identity, and he does not apply the semantics to substitutivity issues.

²⁷ Well, not quite all. We still have various inferences concerning quantifiers, e.g.: $KMc \models \exists x Kx$. We also have those inferences that hold in virtue of the properties of the accessibility relation. Thus, if R is reflexive, $K\alpha \models \alpha$, etc. Some (e.g., Horcutt (1972)), have wondered whether something that verifies so few inferences concerning knowledge is worth calling a *logic* at all. Perhaps not, though even the null logic is, strictly speaking, a logic. More to the point, the fact that the logic is relatively uninteresting does not mean the semantics is uninteresting. It is, in fact, a very hard matter to give an account of the semantics of 'know', which shows why some of the inferences that one might have thought to hold, do not do so – crucially, in the present case, the substitutivity of identicals.

²⁸ And, note, this inference certainly does not work for fear. I can fear meeting Jack the Ripper, and also fear that Jack the Ripper is my wife, without fearing meeting my wife.

²⁹ Even *ideally rational* agents may know that Mn and not know that Mc .

³⁰ In particular, Kripke's puzzle about belief is solved. Let B , M , m , and g be: 'Pierre believes that', 'is a man', 'Mary Anne Evans' and 'George Elliot', respectively. Then it is easy enough to construct models of either kind where $g = m$, BMg and $B\neg Mm$ all hold. Salmon (1995), p. 5, points out that there may well be another rigid designator, d , such

that $d = m = g$ and Pierre has no beliefs about Md at all. It is easy enough to construct models of either kind where, in addition, both $\neg B M d$ and $\neg B \neg M d$ also hold.

³¹ More: suppose that five minutes after Nescio enters the room, someone reliably tells you of Nescio that he was, in fact, born in Megara. Your *de re* knowledge of Nescio would seem to have changed. Yet this cannot be the case if you already knew of Nescio that he was born in Megara.

³² For an account of λ -terms in the context of quantified modal logic, see Fitting and Mendelsohn (1998), Chs. 9, 10.

³³ If *de re* constructions are represented by λ -terms, we would have $\lambda x(K M x)c \wedge \lambda x(\neg K M x)c$ – and the same for n . But this does not convert into a contradiction. λ -conversion will fail in epistemic contexts for exactly the same reason that substitutivity does.

³⁴ See, e.g., Boh (1993). For Burley on the Hooded Man paradox, see p. 40.

³⁵ See, e.g., Quine (1971).

³⁶ A version of this paper was given as the Annual Alice Ambrose Lazerowitz and Thomas Tymoczko Memorial Logic Lecture, Smith College, November 2001. I am grateful to Jay Garfield for making this possible, and to Lee Bowie, who was the commentator on that occasion. Other versions of the paper have been given at the University of Queensland, the University of St Andrews, the Graduate Center, City University of New York, and to a joint meeting of logicians from the Universities of Adelaide and Melbourne. I am grateful to many of those present for their helpful thoughts and comments, and especially to Allen Hazen, Jesper Kallestrup, Arnie Koslow, Calvin Normore, Agustin Rayo, Stephen Read, John Skorupski, Barry Taylor, Achille Varzi, and Crispin Wright; also to an anonymous referee from this journal.

REFERENCES

- Boër, S. E. and Lycan, W. G.: 1986, *Knowing Who*, MIT Press, Cambridge, MA.
- Boh, I.: 1993, *Epistemic Logic in the Later Middle Ages*, Routledge, London.
- Bressan, A.: 1972, *A General Interpreted Modal Calculus*, Yale University Press, New Haven, CT.
- Fitting, M. and Mendelsohn, R.: 1998, *First Order Modal Logic*, Kluwer Academic Publishers, Dordrecht.
- Frege, G.: 1952, On sense and reference, in P. Geach and M. Black (trans.), *Translations from the Writings of Gottlob Frege*, Basil Blackwell, Oxford.
- Gupta, A.: 1980, *The Logic of Common Nouns*, Yale University Press, New Haven, CT.
- Hicks, R. D. (trans.): 1925, *Diogenes Laertius: Lives of Eminent Philosophers*, Harvard University Press, Cambridge, MA.
- Hintikka, J.: 1962, *Knowledge and Belief: An Introduction to the Logic of the Two Notions*, Cornell University Press, Ithaca, NY.
- Hintikka, J.: 1971, Semantics for propositional attitudes, Ch. 10 of Linsky (1971).
- Horcutt, M. O.: 1972, Is epistemic logic possible?, *Notre Dame J. Formal Logic* **13**, 433–453.
- Hughes, G. E. and Cresswell, M.: 1968, *Introduction to Modal Logic*, Methuen, London.
- Kaplan, D.: 1971, Quantifying in, Ch. 9 of Linsky (1971).
- Kneale, W. and Kneale, M.: 1962, *The Development of Logic*, Clarendon Press, Oxford.

- Kripke, S.: 1979, A puzzle about belief, in A. Margalit (ed.), *Meaning and Use*, Reidel, Dordrecht, pp. 239–283.
- Kripke, S.: 1980, *Naming and Necessity*, Basil Blackwell, Oxford.
- Lewis, D.: 1968, Counterpart theory and quantified modal logic, *J. Philos.* **65**, 113–126; reprinted as Ch. 3 of *Philosophical Papers*, Vol. 1, Oxford University Press, Oxford, 1983.
- Linsky, L.: 1971, *Reference and Modality*, Oxford University Press, Oxford.
- Parks, Z.: 1974, Semantics for contingent identity systems, *Notre Dame J. Formal Logic* **15**, 333–334.
- Parks, Z.: 1976, Investigations into quantified modal logic, *Studia Logica* **35**, 109–125.
- Parks, Z. and Smith, T. L.: 1974, The inadequacy of Hughes and Cresswell's semantics for contingent identity systems, *Notre Dame J. Formal Logic* **15**, 331–332.
- Priest, G.: 2000, Objects of thought, *Australasian J. Philos.* **78**, 494–502.
- Priest, G.: 2001, *An Introduction to Non-Classical Logic*, Cambridge University Press, Cambridge.
- Quine, W. V. O.: 1971, Quantifiers and propositional attitudes, Ch. 8 of Linsky (1971).
- Rantala, V.: 1982, Impossible world semantics and logical omniscience, *Acta Philosophica Fennica* **35**, 106–115.
- Salmon, N.: 1986, *Frege's Puzzle*, Ridgeview Publishing Co., Atascadero, CA.
- Salmon, N.: 1995, Being in two minds: Belief with doubt, *Noûs* **29**, 1–20.
- Yi, B.: 1999, Is two a property?, *J. Philos.* **96**, 163–190.

*Department of Philosophy,
University of Melbourne,
Melbourne, Victoria 3010, Australia
e-mail: g.priest@unimelb.edu.au*